

PostgreSQL实战



谭峰 张文升 编著



机械工业出版社
China Machine Press

□□□□□□

PostgreSQL□□

□□ □□□ □□

ISBN□978-7-111-60346-7

□□□□□□□□□□□□2018□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□

□□□□□+ 86-10-68995265

□□□□□service@bbbvip.com

□□□□□www.hzmedia.com.cn

□□□□ @□□□□

□□□□□ □□□□□□□□□□hzebook□

□□

□□

□□

□□□

□1□ □□□□□□

1.1 □□PostgreSQL

1.2 □□PostgreSQL

1.3 □□□□□□□□□□

1.4 □□□□□□□

1.5 □□□□□□□□□□

1.6 □□□□□□□

1.7 □□□□

□2□ □□□□□

2.1 pgAdmin 4□□

2.2 psql□□□□□

2.3 □□□□

□3□ □□□□

3.1 □□□□

3.2 □□□□

3.3 □□/□□□□

3.4 □□□□

3.5 □□□□□□

3.6 □□□□

3.7 □□□□

3.8 json/jsonb□□

- 3.9 数据库索引
- 3.10 数据库安全
- 第4章 SQL数据库
- 4.1 WITH子句
- 4.2 窗口函数
- 4.3 RETURNING子句
- 4.4 UPSERT
- 4.5 数据库连接
- 4.6 数据库备份
- 4.7 数据库迁移
- 4.8 数据库性能优化

第5章 数据库设计

- 第5章 数据库设计
- 5.1 数据库设计概述
- 5.2 数据库设计原则
- 5.3 数据库设计方法
- 5.4 数据库设计工具
- 第6章 数据库应用
- 6.1 数据库应用概述
- 6.2 数据库应用开发
- 6.3 数据库应用部署
- 6.4 数据库应用维护
- 6.5 数据库应用安全

- 第7章 数据库进阶
- 7.1 数据库进阶概述
- 7.2 PostgreSQL数据库

7.3 PostgreSQL

7.4

8

8.1

8.2

8.3

8.4

9 PostgreSQL NoSQL

9.1 jsonb

9.2 json jsonb

9.3 json jsonb

9.4

10

10.1

10.2

10.3

10.4

11 pgbench

11.1

11.2 pgbench

11.3

12

12.1

12.2

12.3

- 12.4 配置
- 12.5 配置
- 12.6 配置
- 12.7 配置
- 12.8 配置
- 12.9 配置
- 12.10 配置
- 12.11 配置
- 13 配置
 - 13.1 配置
 - 13.2 配置
 - 13.3 配置
 - 13.4 SQL配置
 - 13.5 配置
- 14 配置
 - 14.1 Pgpool-II+配置
 - 14.2 Keepalived+配置
 - 14.3 配置
- 15 配置
 - 15.1 配置
 - 15.2 配置
 - 15.3 配置
 - 15.4 配置
- 16 配置
 - 16.1 CREATE EXTENSION
 - 16.2 pg_stat_statements

- 16.3 auto_explain
- 16.4 pg_prewarm
- 16.5 file_fdw
- 16.6 postgres_fdw
- 16.7 Citus
- 16.8 [pg_hba.conf](#)
- 17 Oracle and PostgreSQL
 - 17.1 [pg_hba.conf](#)
 - 17.2 [pg_ident.conf](#)
 - 17.3 [pg_stat_statements](#)
 - 17.4 [pg_stat_statements](#)
 - 17.5 [pg_stat_statements](#)
 - 17.6 [pg_stat_statements](#)
 - 17.7 oracle_fdw [pg_hba.conf](#)
 - 17.8 [pg_hba.conf](#)
- 18 PostGIS
 - 18.1 [pg_hba.conf](#)
 - 18.2 [PostGIS](#)
 - 18.3 [PostGIS](#)
 - 18.4 [PostGIS](#)
 - 18.5 [PostGIS](#)



PostgresPG
PostgreSQL
PostgreSQLPG

PG2011
6PG
PostgreSQL
“Postgres”

2
PostgreSQLPostgreSQL

2011PostgreSQL
PostgreSQLPostgreSQL
7

LeaderBruce MomjianOleg
BartunovSimon RiggsMagnus Hagander

Postgres-XC Leader

Postgres-XL Leader

PGPool Leader

PGStrom Leader.....

PostgreSQL
PGer
PostgreSQL
PostgreSQL

PostgreSQL 10
PostgreSQL
NoSQL

PostgreSQL
PostgreSQL
PostgreSQL

PostgreSQL
PostgreSQL
PostgreSQL

PostgreSQL
PostgreSQL PGER

Postgres

Postgres 2015-2018

2018 4 25



PostgreSQL
PostgreSQL
PostgreSQL
PostgreSQL
PostgreSQL

2010 PostgreSQL DBA
PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL DBA

PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL PostgreSQL
PostgreSQL PostgreSQL



PostgreSQL
14

PostgreSQLデータベースのインストールと設定
SQLデータベースのインストールと設定5
PostgreSQLデータベースのインストールと設定
データベースのインストールと設定10
PostgreSQLデータベースのインストールと設定
データベースのインストールと設定
OracleデータベースPostgreSQLデータベースPostGISデータベース18
データベース

1 PostgreSQLデータベースのインストールと設定
データベースのインストールと設定

2 psqlデータベースのインストールと設定
psqlデータベースのインストールと設定
psqlデータベースのインストールと設定

3 PostgreSQLデータベースのインストールと設定/
データベースのインストールと設定
json/jsonb
データベースのインストールと設定

4 PostgreSQLデータベースのインストールと設定
WITHデータベースのRETURNINGデータベース
UPSERTデータベースのインストールと設定

5 PostgreSQLデータベースのインストールと設定
PostgreSQLデータベースのインストールと設定

6 PostgreSQL 安装与配置

7 PostgreSQL 性能调优

8 PostgreSQL 备份与恢复

9 PostgreSQL 与 NoSQL 对比

10 PostgreSQL 在 Linux 上的部署

11 PostgreSQL 性能测试工具 pgbench

12 PostgreSQL 集群部署

13 PostgreSQL 高可用架构

14 PostgreSQL 与 Pgpool-II 结合使用

□□

□15□□□□PostgreSQL□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□16□□□□□□□□□□□□□□□□□□□□file_fdw□
pg_stat_statements□auto_explain□
postgres_fdw□□□□□□Citus□□□□□

□17□□□□□□□□□□□□□□□□Oracle□□□□□□
PostgreSQL□□□□□□□□□□

□18□□□□□□PostGIS□□□□□□□□□□□□□□□□
□□□□□□□□□□PostGIS□□□□□□□□□□□□□□□□

□□□□

□□□□PostgreSQL□□□□□□□□□□PostgreSQL
□□□□□□□□□□□□□□□□PostgreSQL□□□□□□□□□□
PostgreSQL□□□□□□□□□□□□□□□□□□□□□□
PostgreSQL□□□□□□□□□□□□

□□□□PostgreSQL 10□□□□□□□□□□□□
PostgreSQL 10□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□
PostgreSQL 10□□□□□□□□□□

18 PostgreSQL
PostgreSQL
PostgreSQL
PostgreSQL

PostgreSQL

PostgreSQL DBA PostgreSQL DBA PostgreSQL

PostgreSQL DBA Oracle MySQL DBA PostgreSQL PostgreSQL

PostgreSQL PostgreSQL PostgreSQL PostgreSQL

PostgreSQL PostgreSQL PostgreSQL

--	--	--	--	--

francs.tan@postgres.cn

11

A diagram showing a 4x28 grid of squares. The top row has 28 squares. The second row has 28 squares. The third row has 28 squares. The fourth row has 10 squares, starting from the left.

□□□francs□

2018□4□□□□

□□□

・□1□ □□□□□□□

・□2□ □□□□□

・□3□ □□□□

・□4□ SQL□□□□

1 数据库

PostgreSQL数据库

1.1 PostgreSQL

PostgreSQL

- UC Berkeley 1977
- Ingres Michael Stonebraker 1994
- UC Berkeley Andrew Yu Jolly Chen
- SQL Ingres QUEL
- Postgres95 SQL
- Postgres95 1996 PostgreSQL
- PostgreSQL 6.0 2005
- PostgreSQL Windows 8.0
- 2010 PostgreSQL 9.0 PostgreSQL
- PostgreSQL
- PostgreSQL 10

PostgreSQL

PostgreSQL RDS



 ☐ PostgreSQL ☐ "Post-Gres-Q-L" ☐ `[/postgrs kju|/]` ☐ PostgreSQL ☐ Postgres

PostgreSQL
<http://www.postgresql.org/files/postgresql.mp3>

1.1.1 PostgreSQL

PostgreSQL Linux
UNIX UNIX Windows AIX
BSD HP-UX SGI IRIX Mac OS X Solaris
Tru64 C C++ Go Java Perl
Python Ruby Tcl ODBC

json jsonb
SQL CREATE TYPE

SQL SQL SQL
Window Function
C Java python R

MVCC
Hot Standby
Foreign data wrappers
Oracle MySQL Informix
SQLite MS SQL Server

1.1.2 许可

PostgreSQL 采用 PostgreSQL License 许可证。该许可证是宽松的，允许用户自由使用、复制、分发和修改 PostgreSQL，只要保留原始的版权声明和许可证副本即可。该许可证与 BSD、MIT 和 OSI 许可证类似。PostgreSQL 的源代码和文档都遵循该许可证。PostgreSQL 的文档和源代码都遵循该许可证。PostgreSQL 的文档和源代码都遵循该许可证。

有关 PostgreSQL 许可证的更多信息，请访问：
<https://www.postgresql.org/about/licence/>

1.1.3 社区资源

PostgreSQL 有一个活跃的社区，提供许多资源。包括 PostgreSQL 的官方文档、社区论坛、邮件列表、书籍、教程、视频等。PostgreSQL 的社区资源非常丰富，可以帮助用户学习和使用 PostgreSQL。PostgreSQL 的社区资源非常丰富，可以帮助用户学习和使用 PostgreSQL。

PostgreSQL 的官方文档和教程，请访问：
<https://www.postgresql.org/list>

1.2 PostgreSQL

PostgreSQLは、オープンソースのデータベースシステムです。32ビット/64ビットのアーキテクチャをサポートし、CentOS 6にPostgreSQL 10をインストールします。

PostgreSQLは、yumを使用してインストールできます。インストールコマンドは以下のとおりです。

1.2.1 yum

yumを使用してPostgreSQLをインストールするには、以下のコマンドを実行します。

```
sudo yum install postgresql
```

インストールが完了すると、PostgreSQLのサービスが自動的に起動されます。

1. PostgreSQL repository RPM

PostgreSQLのRPMパッケージは、<https://www.postgresql.org/download> からダウンロードできます。CentOSのインストールには、"Binary packages"セクションの"Red Hat family Linux"を選択し、バージョンとして"10"を選択し、プラットフォームとして"CentOS 6"を選択します。

architecture="x86_64" repository RPM

```
[root@pghost1 ~]$ yum install  
https://download.postgresql.org/pub/repos/yum/10/redhat/rhel  
-6-x86_64/pgdg-centos10-10-1.noarch.rpm
```

/etc/yum.repos.d/pgdg-10-centos.repo

2. PostgreSQL

PostgreSQL repository RPM
yum search postgresql10

```
[root@pghost1 ~]$ yum search postgresql10
```

·postgresql10-debuginfo.x86_64
postgresql10 DEBBUG

·postgresql10.x86_64
PostgreSQL client

·postgresql10-contrib.x86_64
PostgreSQL 라이브러리


·postgresql10-devel.x86_64
PostgreSQL C/C++ 라이브러리 libpq

·postgresql10-docs.x86_64

·postgresql10-server.x86_64
PostgreSQL server

server contrib
client

```
[root@pghost1 ~]$ yum install postgresql10-server  
postgresql10-contrib
```

 `yum -y`

```
[root@pghost1 ~]$ yum install -y postgresql10-server  
postgresql10-contrib
```

yum
/usr/pgsql-10
/usr/pgsql-10/bin
postgres
home
/var/lib/pgsql

3. yum PostgreSQL

PostgreSQL

```
[root@pghost1 ~]$ rpm -qa | grep postgresql  
postgresql10-10.0-1PGDG.rhel6.x86_64  
postgresql10-libs-10.0-1PGDG.rhel6.x86_64  
postgresql10-contrib-10.0-1PGDG.rhel6.x86_64  
postgresql10-server-10.0-1PGDG.rhel6.x86_64
```

yum remove
libs
libs

```
yum remove postgresql10-libs-10.0-1PGDG.rhel6.x86_64
```

PostgreSQL
init.d

```
[root@pghost1 ~]$ rm -f /etc/init.d/postgresql-10
```

1.2.2

如何安装PostgreSQL数据库

1. 安装

PostgreSQL官网
<https://www.postgresql.org/ftp/latest> 下载
postgresql-10.0.tar.gz

```
[root@pghost1 ~]$ wget  
https://ftp.postgresql.org/pub/source/v10.0/postgresql-  
10.0.tar.gz
```

解压

```
[root@pghost1 ~]$ tar -xvf postgresql-10.0.tar.gz
```

2. 安装configure

安装configure依赖包

```
[root@pghost1 ~]$ yum groupinstall "Development tools"  
[root@pghost1 ~]$ yum install -y bison flex readline-devel  
zlib-devel
```

configure--help
[root@pghost1 ~]\$ cd postgresql-10.0
[root@pghost1 ~]\$./configure --help | less

PostgreSQL

·--prefix=PREFIX
"/usr/local/pgsql"

·--includedir=DIR
"PREFIX/include"

·--with-pgport=PORTNUM
RPM

·--with-blocksize=BLOCKSIZE
8kB OLAP
32kB OLAP OLTP 8kB

·--with-segsize=SEGSIZE
1GB

found”yum

3.

Linux PostgreSQL GNU
make gmake gmake install
gmake world gmake install-
world world
world

gmake gmake world

```
[root@pghost1 ~]$ gmake
```

gmake “All of
PostgreSQL successfully made. Ready to
install.”

gmake world
“PostgreSQL contrib and documentation
successfully made. Ready to install.”

```
gmake installgmake install-world
```

```
[root@pghost1 ~]$ gmake install
```

```
gmake install
"PostgreSQL installation complete."
```

```
gmake install-world
"PostgreSQL contrib and
documentation installation complete."
```

```
PostgreSQL
```

```
[root@pghost1 ~]# /opt/pg10/bin/postgres --version
postgres (PostgreSQL) 10.0
```

1.2.3

```
shellPython
/opt/pg9.x

```

ln -s /opt/pgsql

```
[root@pghost1 opt]$ ln -s /opt/pg10 /opt/pgsql
[root@pghost1 ~]$ ll /opt/
drwxr-xr-x  6 root    root    4096 Oct 11 14:32 pg96
drwxr-xr-x  6 root    root    4096 Oct 13 17:43 pg10
lrwxrwxrwx  1 root    root     10 Oct 13 11:25 pgsql ->
/opt/pg10/
```

ln -s /opt/pg10 /opt/pgsql

1.3 安装PostgreSQL

安装PostgreSQL

```
[postgres@pghost1 ~]$ tree -L 1 /opt/pgsql/  
/opt/pgsql/  
├── bin  
├── include  
├── lib  
└── share  
4 directories, 0 files
```

share 包含 PostgreSQL 的 man 手册
include 包含 PostgreSQL 的 C/C++ 头文件
bin 包含 PostgreSQL 的可执行文件
lib 包含 PostgreSQL 的库文件
C/S 架构的数据库系统，由数据库引擎和数据库客户端组成。
数据库引擎负责存储和管理数据，数据库客户端负责与引擎交互。
数据库引擎和数据库客户端之间通过通信协议进行交互。

1.3.1 安装

安装PostgreSQL

1. 安装SQL数据库

clusterdb

clusterdb SQL CLUSTER
PostgreSQL clusterdb
clusterdb SQL CLUSTER

clusterdb

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/clusterdb -h pghost1 -p 1921 -d mydb
```

reindexdb

reindexdb SQL REINDEX
reindexdb SQL REINDEX
reindexdb SQL REINDEX

reindexdb

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/reindexdb -e -h pghost1 -p 1921 -d mydb
```

vacuumdb

vacuumdb PostgreSQL
VACUUM VACUUM FREEZE VACUUM
FULL VACUUM ANALYZE SQL

VACUUM PostgreSQL

vacuumdb

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/vacuumdb -h pghost1 -p 1921 mydb
```

vacuumlo

vacuumlo

vacuumlo

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/vacuumlo -h pghost1 -p 1921 mydb
```

createdb dropdb

SQL CREATE DATABASE DROP DATABASE

pghost1 1921 newdb

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/createdb -h pghost1 -p 1921 newdb "New database."
```

newdb

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/dropdb -h pghost1 -p 1921 newdb
```

createuser dropuser

SQL CREATE USER DROP USER

newuser newuser
pg_monitor 1

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/createuser -h pghost1 -p 1921 -c 1 -g pg_monitor -D -R -S -P -e newuser
Enter password for new role:
Enter it again:
CREATE ROLE newuser PASSWORD
'md518b2c3ec6fb3de0e33f5612ed3998fa4' NOSUPERUSER NOCREATEDB
NOCREATEROLE INHERIT LOGIN CONNECTION LIMIT 1 IN ROLE
pg_monitor;
```

--interactive

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/createuser -h pghost1 -p 1921 -c 1 -g pg_monitor --interactive -e -P newuser
```

```
dropuser newuser
```

3. 問題点

ecpg C PostgreSQL SQL
SQL C SQL C
C C

·oid2name PostgreSQL OID
OID

·pgbench
pgbench

·pg_config PostgreSQL
pg_config

·PostgreSQL pg_isready
pg_isready

·pg_receivexlog
pg_receivexlog

·pg_recvlogical
pg_recvlogical

·psql PostgreSQL
psql psql
psql

1.3.2 配置

配置

·initdb

·pg_archivecleanup PostgreSQL
WAL

·pg_controldata
pg_control

·pg_ctl

·pg_resetwal
pg_control
pg_resetwal

·pg_rewind master slave
master

·pg_test_fsync
wal_sync_method
I/O

·pg_test_timing

·pg_upgrade PostgreSQL 10.5 10.6
10.5 10.6

·pg_waldump 10.5 10.6 10.7 10.8 10.9 11.0

·postgres PostgreSQL 10.5 10.6

·postmaster 10.5 bin 10.6 10.7 10.8 10.9 postgres
10.5 10.6 10.7 10.8 10.9

1.4 安装PostgreSQL

PostgreSQL 是一个开源的数据库系统，它支持多种数据类型，并且具有强大的查询能力。在安装 PostgreSQL 之前，需要先安装一些依赖包。在 Ubuntu 系统中，可以使用以下命令来安装 PostgreSQL 及其依赖包：

1.4.1 安装依赖包

在安装 PostgreSQL 之前，需要先安装一些依赖包。在 Ubuntu 系统中，可以使用以下命令来安装 PostgreSQL 及其依赖包：

```
sudo apt-get install postgresql postgresql-contrib
```

安装完成后，可以使用以下命令来启动 PostgreSQL 服务：

```
sudo systemctl start postgresql
```

如果服务已经启动，可以使用以下命令来检查服务的状态：

```
sudo systemctl status postgresql
```

如果服务没有启动，可以使用以下命令来启动服务：

```
sudo systemctl start postgresql
```

如果服务已经启动，可以使用以下命令来检查服务的状态：

```
sudo systemctl status postgresql
```

如果服务没有启动，可以使用以下命令来启动服务：

```
sudo systemctl start postgresql
```

在安装 PostgreSQL 之前，需要先安装一些依赖包。在 Ubuntu 系统中，可以使用以下命令来安装 PostgreSQL 及其依赖包：

```
sudo apt-get install postgresql postgresql-contrib
```

安装完成后，可以使用以下命令来启动 PostgreSQL 服务：

```
sudo systemctl start postgresql
```

如果服务已经启动，可以使用以下命令来检查服务的状态：

```
sudo systemctl status postgresql
```

如果服务没有启动，可以使用以下命令来启动服务：

```
sudo systemctl start postgresql
```

PostgreSQLのBUGについて

インストール

```
[root@pghost1 ~]$ groupadd -g 1000 postgres
[root@pghost1 ~]$ useradd -g 1000 -u 1000 postgres
[root@pghost1 ~]$ id postgres
uid=1000(postgres) gid=1000(postgres) groups=1000(postgres)
```

確認

1 rootでsudoを実行

2 NTPのインストール

```
uid gid groupmod usermod
```

```
[root@pghost1 ~]$ groupmod -g 1000 postgres
[root@pghost1 ~]$ usermod -u 1000 -g 1000 postgres
```

1.4.2 インストール

Diagram illustrating a grid of 50 empty rectangular boxes (5 rows by 10 columns). The second box in the second row contains the text "SQL".

```

pgdata/9.x/xxx_data
pgdata/9.x/10/xxx_data
pgdata/10/data
backups
scripts
archive_wals

```

```
[root@pghost1 ~]$ mkdir -p
/pgdata/10/{data,backups,scripts,archive wals}
```

```

0700 initdb PostgreSQL

```

```
[root@pghost1 ~]$ chown -R postgres.postgres /pgdata/10
[root@pghost1 ~]$ chmod 0700 /pgdata/10/data
```

1.4.3 安装数据库

安装数据库使用initdb命令。initdb命令的默认配置在/usr/share/doc/postgresql/template1/postgresql.conf文件中。initdb命令的默认配置在/usr/share/doc/postgresql/template1/postgresql.conf文件中。

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/initdb --help
initdb initializes a PostgreSQL database cluster.
Usage:
    initdb [OPTION]... [DATADIR]
Options:
    -A, --auth=METHOD          认证方法，默认为md5
                                md5
                                trust
                                trust
                                认证方法，默认为md5
                                认证方法，默认为md5
                                认证方法，默认为md5
    --auth-host=METHOD         认证方法，默认为md5
                                认证方法，默认为md5
    pg_hba.conf                认证方法，默认为md5
    --auth-local=METHOD        认证方法，默认为md5
    pg_hba.conf                认证方法，默认为md5
    [-D, --pgdata=]DATADIR     数据目录
                                数据目录
    -E, --encoding=ENCODING     字符集
                                字符集
                                字符集
    --locale=LOCALE             本地化设置
                                本地化设置
    --lc-collate=, --lc-ctype=, --lc-messages=LOCALE
    --lc-monetary=, --lc-numeric=, --lc-time=LOCALE
                                本地化设置
    --no-locale                 本地化设置
                                本地化设置
    --pwfile=FILE               密码文件
                                密码文件
    -T, --text-search-config=CFG 全文搜索配置
                                全文搜索配置
    -U, --username=NAME         用户名
                                用户名
    -W, --pwprompt               密码提示
                                密码提示
    -X, --waldir=WALDIR         WAL目录
```

template1 postgres template1
postgres
postgres
postgres

postgres
postgres

initdb: directory "/pgdata/10/data" exists but is not empty
If you want to create a new database system, either remove
or empty
the directory "/pgdata/10/data" or run initdb
with an argument other than "/pgdata/10/data".

postgres

initdb postgresql_ctl
postgres

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_ctl init -D  
/pgdata/10/data -o "-W"
```

postgres

yum PostgreSQL
/var/lib/pgsql/10 data
backups service postgresql-10 init

/var/lib/pgsql/10/data

1.5 PostgreSQL의 서비스

PostgreSQL을 시스템 서비스로 등록하고, service를 사용하여 PostgreSQL을 시작하고, 중지하는 방법을 소개합니다.

1.5.1 service를 사용하여 PostgreSQL을 시작하고 중지하는 방법

PostgreSQL을 시스템 서비스로 등록합니다.

```
[root@pghost1 ~]$ service postgresql-10 start
```

PostgreSQL을 시스템 서비스로 등록합니다.

```
[root@pghost1 ~]$ service postgresql-10 status
```

PostgreSQL을 시스템 서비스로 등록합니다.

```
[root@pghost1 ~]$ service postgresql-10 stop
```

1.5.2 pg_ctl을 사용하여 PostgreSQL을 시작하고 중지하는 방법

pg_ctl PostgreSQL 数据库管理系统
service
systemctl pg_ctl
pg_ctl su postgres

1. 启动

启动

```
[root@pghost1 ~]# su - postgres
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_ctl -D
/pgdata/10/data start
server started
```

2. 检查

检查

```
[root@pghost1 ~]# su - postgres
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_ctl -D
/pgdata/10/data status
pg_ctl: no server running
```

检查

```
pg_ctl: server is running (PID: 43965)
/opt/pgsql/bin/postgres "-D" "/pgdata/10/data"
```

pg_isready 检查数据库是否接受连接

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_isready -p 1921  
/tmp:1921 - accepting connections
```

pg_isready

```
/tmp:1921 - no response
```

3. pg_ctl

pg_ctl 管理 PostgreSQL 数据库

```
pg_ctl stop      [-D DATADIR] [-m SHUTDOWN-MODE] [-W] [-t  
SECS] [-s]
```

“-s” 指定在 shutdown 时使用的模式，默认为 “-t SECS” 指定在 shutdown 时等待的时间，默认为 SECS 秒。 “-m” 指定在 shutdown 时使用的模式，默认为 smart。 PostgreSQL 支持三种 shutdown 模式： smart、fast、immediate。 fast 模式会快速关闭数据库，但可能会丢失未提交的事务。 immediate 模式会立即关闭数据库，但可能会丢失未提交的事务。

·smart 模式会等待所有事务完成后再关闭数据库，这是默认模式。

·fast 快速启动，启动时只加载共享库，不加载数据文件，启动速度快，但数据文件未加载，数据库无法使用。

·immediate 立即启动，启动时只加载共享库，不加载数据文件，启动速度快，但数据文件未加载，数据库无法使用。

启动时只加载共享库，不加载数据文件，启动速度快，但数据文件未加载，数据库无法使用。
“-ms”“-mf”“-mi” 启动时只加载共享库，不加载数据文件，启动速度快，但数据文件未加载，数据库无法使用。

```
[root@pghost1 ~]# su - postgres
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_ctl -D
/pgdata/10/data -ms stop
```

1.5.3 启动 PostgreSQL 数据库

启动 PostgreSQL 数据库，需要启动 postmaster 进程，postgres 进程。

```
[root@pghost1 ~]# su - postgres
[postgres@pghost1 ~]$ /opt/pgsql/bin/postgres -D
/pgdata/10/data/
```

启动 PostgreSQL 数据库，需要启动 postmaster 进程，postgres 进程。

PostgreSQL 数据库启动后，postmaster 进程会接收 SIGINT、SIGTERM、SIGQUIT 信号。

PostgreSQL smart fast
immediate kill postgres
SIGTERM SIGINT SIGQUIT
smart

```
[postgres@pghost1 ~]$ kill -sigterm `head -1  
/pgdata/10/data/postmaster.pid`  
received smart shutdown request  
shutting down  
database system is shut down
```

smart
PostgreSQL
pqsignal pmdie

PostgreSQL pg_ctl
kill

1.5.4

yum

1.

contrib Linux FreeBSD OSX

```
[root@pghost1 ~]$ ls postgresql-10.0/contrib/start-scripts/  
freebsd  linux  osx
```

将linux下的/etc/init.d/postgresql-10复制到/etc/init.d/

```
[root@pghost1 ~]$ cp postgresql-10.0/contrib/start-  
scripts/linux /etc/init.d/postgresql-10  
[root@pghost1 ~]$ chmod +x /etc/init.d/postgresql-10  
[root@pghost1 ~]$ ls -lh /etc/init.d/postgresql-10  
-rwxr-xr-x 1 root root 3.5K Oct 13 16:30  
/etc/init.d/postgresql-10
```

2. 检查配置

chkconfig--list查看PostgreSQL配置

```
[root@pghost1 ~]$ chkconfig --list | grep postgresql-10  
postgresql-10 0:off 1:off 2:off 3:off 4:off  
5:off 6:off
```

chkconfig查看PostgreSQL配置

```
[root@pghost1 ~]$ chkconfig postgresql-10 on/off
```

1.6 安装配置

安装配置 PostgreSQL 数据库系统，首先需要安装 PostgreSQL 数据库系统，然后配置 Database 系统，最后配置 PostgreSQL 数据库系统。

PostgreSQL 数据库系统安装配置，首先需要安装 PostgreSQL 数据库系统，然后配置 Database 系统，最后配置 PostgreSQL 数据库系统。安装 PostgreSQL 数据库系统，需要安装 postgresql.conf、pg_hba.conf 等配置文件。配置 Database 系统，需要配置 postgresql.conf、pg_hba.conf 等配置文件。配置 PostgreSQL 数据库系统，需要配置 postgresql.conf、pg_hba.conf 等配置文件。

1.6.1 安装配置

安装配置 PostgreSQL 数据库系统，首先需要安装 PostgreSQL 数据库系统，然后配置 Database 系统，最后配置 PostgreSQL 数据库系统。安装 PostgreSQL 数据库系统，需要安装 postgresql.conf、postgresql.auto.conf、pg_hba.conf、pg_ident.conf 等配置文件。配置 Database 系统，需要配置 postgresql.conf、pg_hba.conf 等配置文件。配置 PostgreSQL 数据库系统，需要配置 postgresql.conf、pg_hba.conf 等配置文件。

1.6.2 pg_hba.conf

pg_hba.conf 配置文件用于配置 PostgreSQL 数据库系统的访问控制。该文件位于 PostgreSQL 数据库系统的安装目录下。该文件用于配置 PostgreSQL 数据库系统的访问控制。

TYPE	DATABASE	USER	ADDRESS
METHOD			
local	database	user	auth-method [auth-options]
host	database	user	address auth-method [auth-options]
hostssl	database	user	address auth-method [auth-options]
hostnossl	database	user	address auth-method [auth-options]
host	database	user	IP-address IP-mask auth-method [auth-options]
hostssl	database	user	IP-address IP-mask auth-method [auth-options]
hostnossl	database	user	IP-address IP-mask auth-method [auth-options]

```

# TYPE DATABASE USER ADDRESS METHOD
# local database user auth-method [auth-options]
# host database user address auth-method [auth-options]
# hostssl database user address auth-method [auth-options]
# hostnossl database user address auth-method [auth-options]
# host database user IP-address IP-mask auth-method [auth-options]
# hostssl database user IP-address IP-mask auth-method [auth-options]
# hostnossl database user IP-address IP-mask auth-method [auth-options]

```

1. 設定

TYPEはlocal、host、hostssl、hostnosslの4種類がある。

・localはUnixのlocalhostに接続するための設定。

・hostはTCP/IPで接続するための設定。SSLはSSLの有無を指定する。localhostはlocalhostに接続するための設定。TCP/IPで接続するための設定はpostgresql.confのlisten_addressesで指定する。

·hostssl□□□□□□□□SSL□TCP/IP□□□□□□
hostssl□□□□□□□□

1.□□□□□□□□□□OpenSSL□

2.□□PostgreSQL□□□□□configure□□--
with-openssl□□SSL□□□□

3.□postgresql.conf□□□ssl=on□

·hostnossll□hostssl□□□□□□□□□□SSL□
TCP/IP□□□□

2.□□□□□□

DATABASE□□□□□□□□□□□□□□□□□□

3.□□□□□□

USER□□□□□□□□□□□□□□□□□□□□

4.□□□□□□

ADDRESS□□□□□□□□□□□□IP□□□□IP□□□□□□□□

5.□□□□□□

METHOD trust
reject md5 password

reject
password

md5 password md5
md5 password
password

scram-sha-256 PostgreSQL 10
SASL PostgreSQL
scram-sha-256
PostgreSQL 10

```
[postgres@pghost2 ~]$ /usr/pgsql-9.6/bin/psql -h pghost1 -p 1921 -U postgres mydb  
psql: SCRAM authentication requires libpq version 10 or above
```

<https://www.postgresql.org/docs/current/static/auth-methods.html>

1.6.3 postgresql.conf

postgresql.conf 파일의 값을 변경하려면
configparameter=value로 값을 설정하고 #로 주석 처리
한 후 value를 변경하고 MB, GB, ms, min, d 등의
value를 지정하여 MB, GB, ms, min, d 등의
postgresql.conf 파일을 include
include_if_exists로 포함시킵니다.

postgresql.conf 파일을 #로 변경하려면
restart로 서버를 재시작하거나
reload로 서버를 재로드합니다.

1. postgresql.conf 파일

postgresql.conf 파일

- postgresql.conf 파일을
- vim, nano, sed 등으로
- ALTER SYSTEM로

```
mydb=# ALTER SYSTEM SET listen_addresses = '*';
```

ALTER SYSTEM SQL 문을 사용하여
postgresql.auto.conf 파일을

postgresql.auto.conf

postgresql.conf

-

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/postgres -D  
/pgdata/10/data -c port=1922
```

2.

- Database

```
ALTER DATABASE name SET configparameter { TO | = } { value |  
DEFAULT }  
ALTER DATABASE name RESET configuration
```

- Session

- SET Session

```
SET configparameter { TO | = } { value | 'value' | DEFAULT }  
SET configparameter TO DEFAULT;
```

- pg_settings

```
UPDATE `pg_settings` SET setting = new_value WHERE name =  
'configparameter';
```

```
UPDATE `pg_settings` SET setting = reset_val WHERE name =  
'configparameter';
```

· `set_config()`

```
SELECT set_config('configparameter',new_value,false);
```

· `ALTER ROLE`

```
ALTER ROLE name IN DATABASE database_name SET  
configparameter { TO | = } { value | DEFAULT }  
ALTER ROLE name IN DATABASE database_name RESET  
configparameter
```

3. `pg_settings`

`pg_settings`

```
SELECT name,setting FROM pg_settings where name ~ 'xxx';  
SELECT current_setting(name);
```

`show` `show all`

4. `reload`

`reload`

```
mydb=# SELECT pg_reload_conf();
pg_reload_conf
-----
t
(1 row)
```

pg_ctl reload

```
[root@pghost1 ~]# su - postgres
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_ctl -D
/pgdata/10/data reload
```

1.6.4

PostgreSQL

```
[postgres@pghost1 ~]$ netstat -nlt | grep 1921
Active Internet connections (only servers)
Proto Recv-Q Send-Q Local Address           Foreign Address         State
tcp        0      0 127.0.0.1:1921          0.0.0.0:*               LISTEN
tcp        0      0 :::1:1921               :::*                     LISTEN
```

```
[postgres@pghost2 ~]$ telnet pghost1 1921
Trying pghost1...
telnet: connect to address pghost1: Connection refused
```



```
·comma-separated list of addresses—  
□□□□□□□□□□
```

```
·defaults to 'localhost' use '*' for all—
"localhost" "*"
VIP "*"

```

- change requires restart —

```
listen_addresses="#"
"*"

```

```
listen_addresses = '*'
```

[illegible]

```
[root@pgghost1 ~]$ service postgresql-10 restart
```

```
pg_hba.conf
pg_hba.conf
TCP/IP
SSL
SSL
md5
pguser
mydb
pg_hba.conf
```

```
[postgres@pghost1 ~]$ echo "host mydb pguser 0.0.0.0/0 md5"
>> /pgdata/10/data/pg_hba.conf
```

pg_hba.conf reload

```
[postgres@pghost1 ~]$ /opt/pgsql/bin/pg_ctl -D
/pgdata/10/data/ reload
server signaled
2017-10-18 10:16:00.405 CST [36171] LOG:  received SIGHUP,
reloading configuration files
```

Windows Linux selinux
iptables Linux selinux
iptables iptables
iptables

1.7 □□□□

[illegible]

2 安装

安装 PostgreSQL 数据库，需要安装 pgAdmin 和 psql。pgAdmin 是一个图形化的 PostgreSQL 数据库管理工具，而 psql 是一个命令行工具。安装 PostgreSQL 数据库，需要安装 PostgreSQL 数据库服务器和 PostgreSQL 客户端工具。安装 PostgreSQL 数据库服务器，需要安装 PostgreSQL 数据库服务器和 PostgreSQL 客户端工具。安装 PostgreSQL 客户端工具，需要安装 PostgreSQL 数据库服务器和 PostgreSQL 客户端工具。

2.1 pgAdmin 4

pgAdmin은 PostgreSQL를 관리하는 도구입니다.
다운로드: <https://www.pgadmin.org/>
pgAdmin 4는 PostgreSQL을 관리하는 데 사용됩니다.

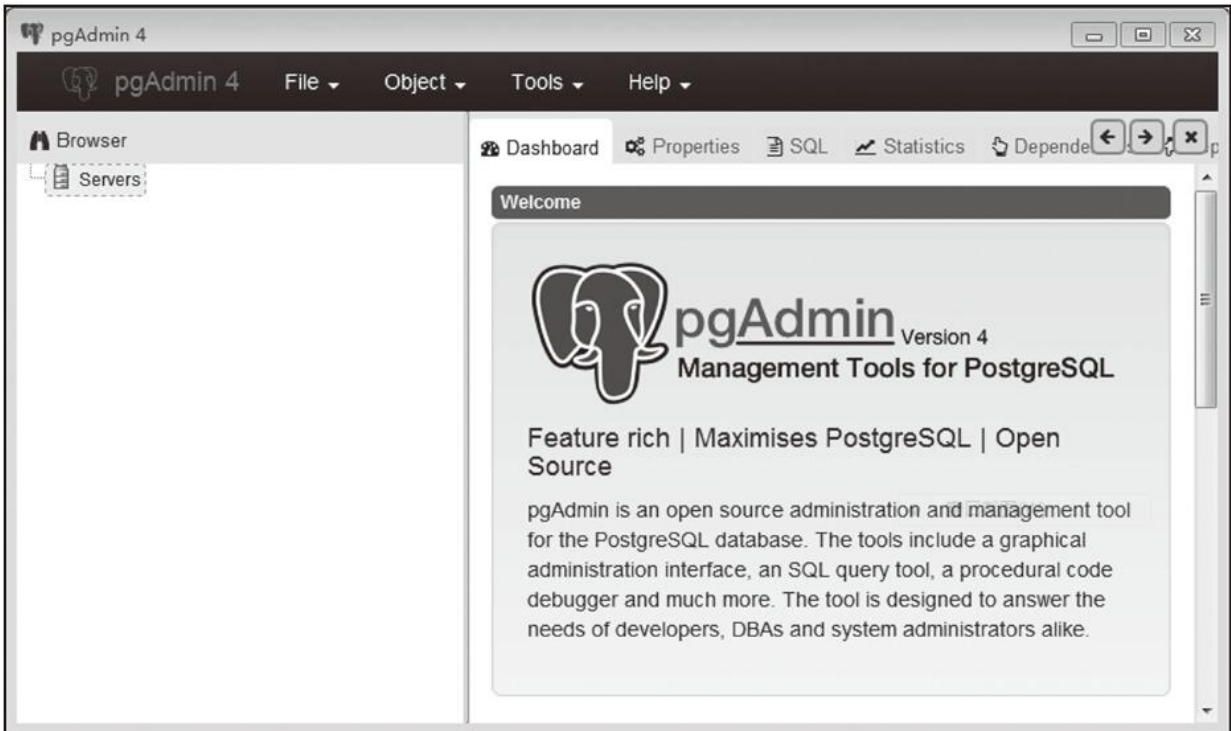
2.1.1 pgAdmin 4

pgAdmin은 Linux, Unix, Mac OS X, Windows에서 실행됩니다. pgAdmin 4는 PostgreSQL을 관리하는 데 사용됩니다. pgAdmin 4는 PostgreSQL을 관리하는 데 사용됩니다. Windows 7에서 pgAdmin 4를 설치하려면 pgAdmin 4

를 다운로드합니다.
<https://www.postgresql.org/ftp/pgadmin/pgadmin4/v1.6/windows/>에서 다운로드합니다.
pgAdmin 4 2-1

2.1.2 pgAdmin 4

pgAdmin 4는 PostgreSQL을 관리하는 데 사용됩니다. pgAdmin 4는 PostgreSQL을 관리하는 데 사용됩니다.



2-1 pgAdmin 4

1.pgAdmin 4

pgAdmin 4 2-2

Create - Server

General **Connection** Advanced

Host name/address 192.168.28.74

Port 1921

Maintenance database postgres

Username postgres

Password

Save password? ☐ 矩形截图(R)

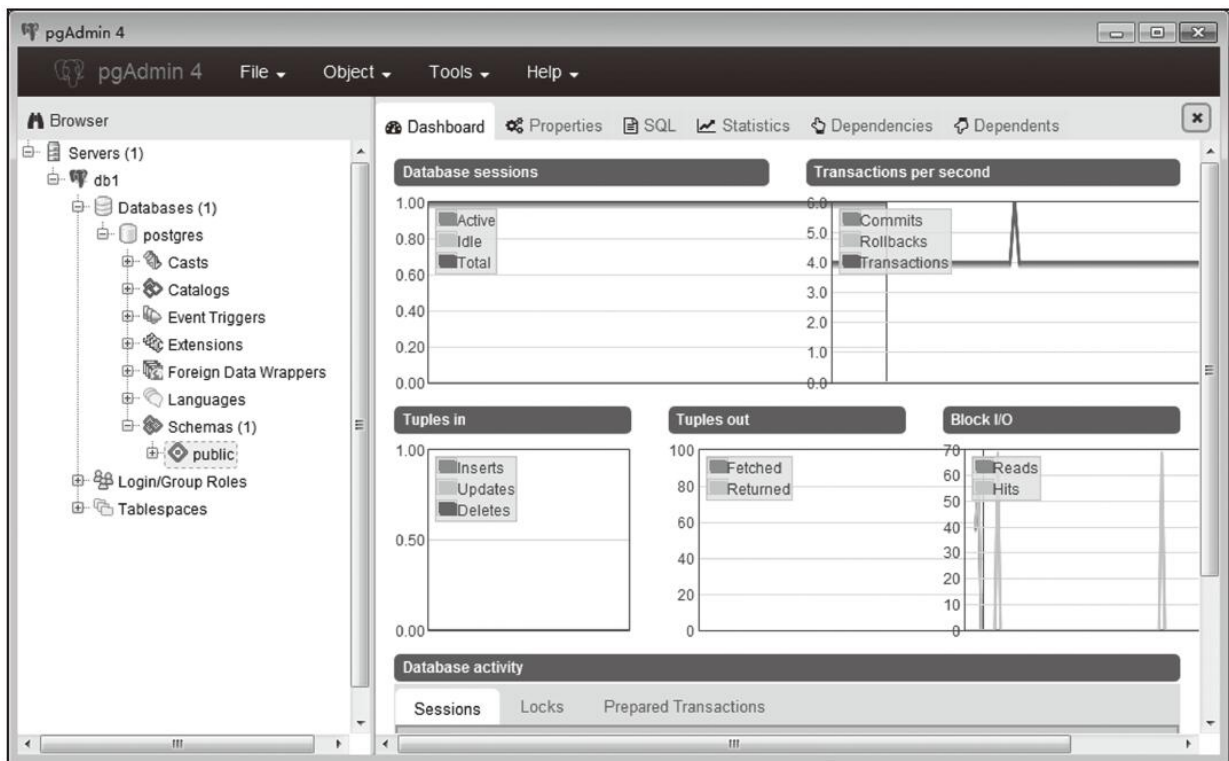
Role

SSL mode Prefer

i ? Save Cancel Reset

图2-2 pgAdmin 4创建数据库

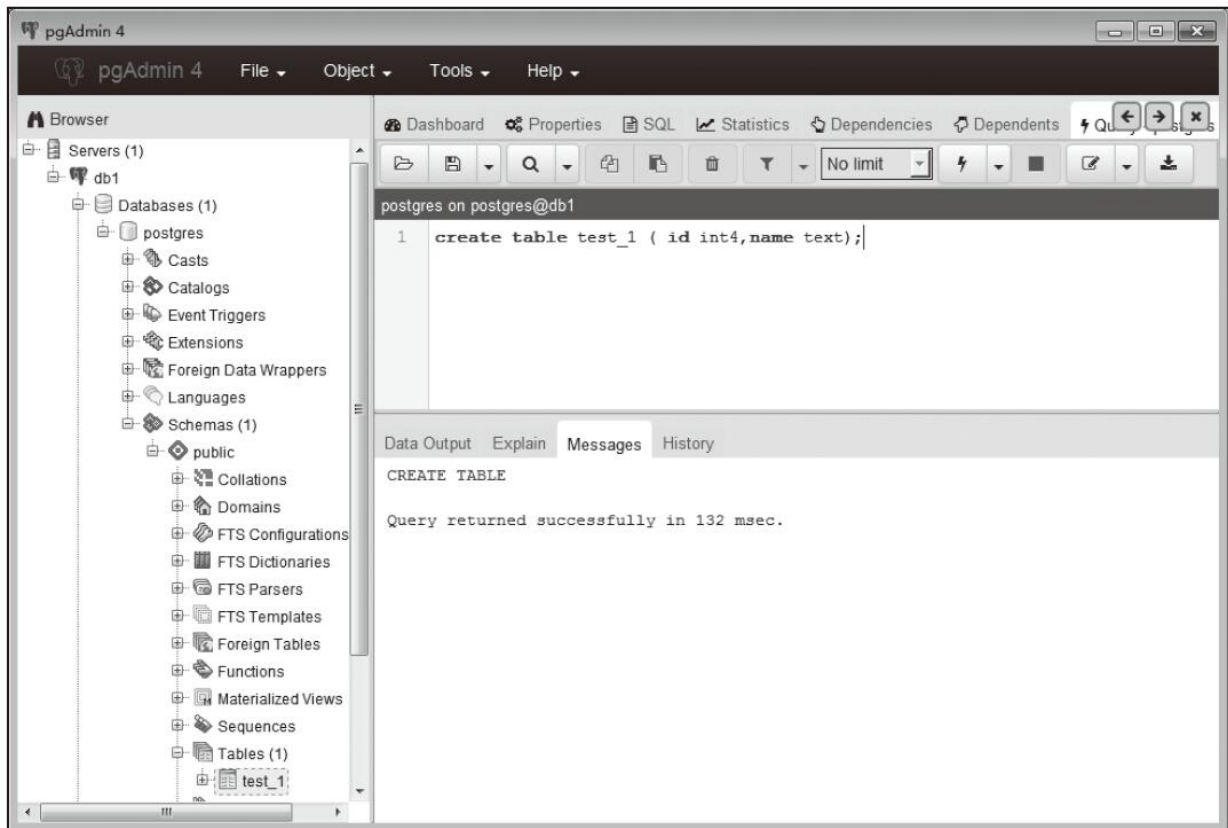
在General选项卡中，将数据库名称设置为db1，
在Connection选项卡中，按照图2-3所示



2-3 pgAdmin 4

2.pgAdmin 4

pgAdmin 4 Tools Query Tool
DDL DML
2-4

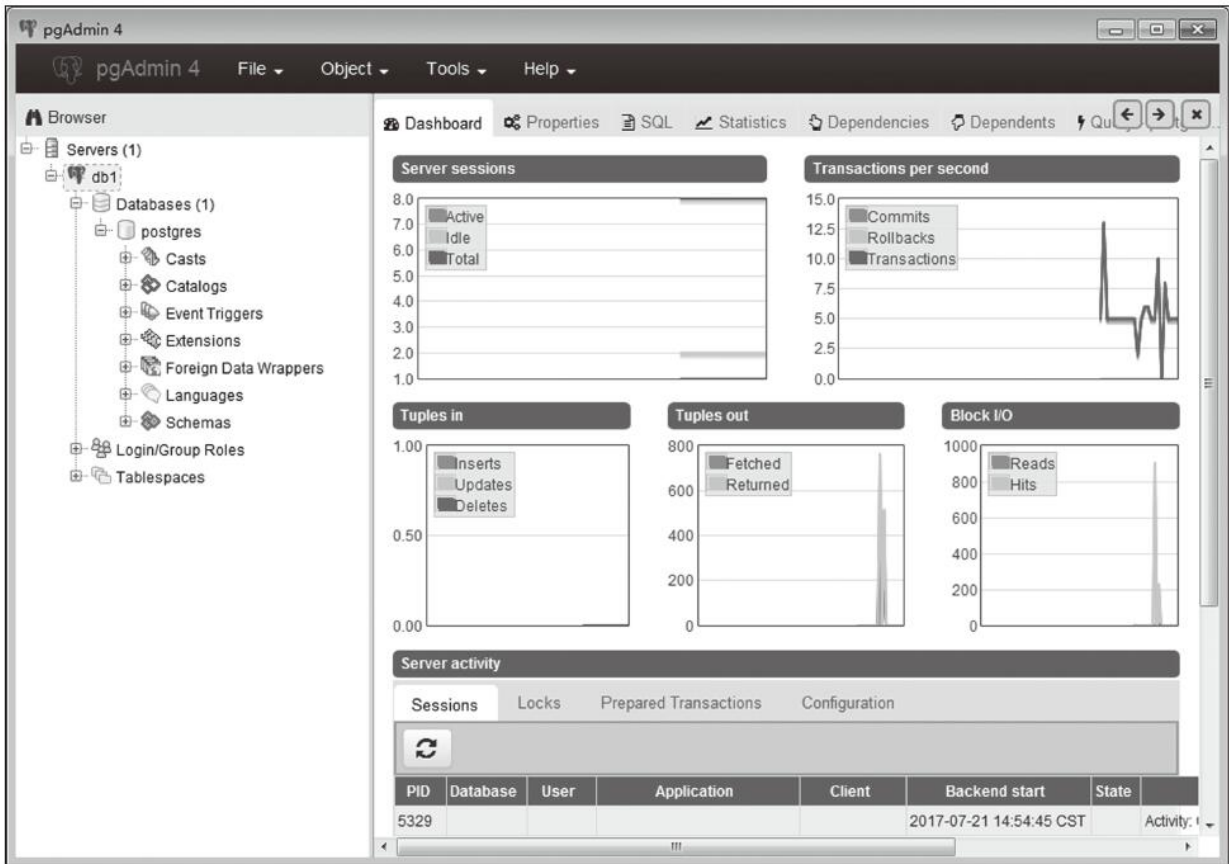


2-4 pgAdmin 4

pgAdmin 4は、PostgreSQLの管理を行うためのツールです。このツールを使用して、データベースの構造を定義（DDL）することができます。

3. pgAdmin 4

pgAdmin 4は、PostgreSQLの管理を行うためのツールです。このツールを使用して、データベースの構造を定義（DDL）することができます。



2-5 pgAdmin 4

pgAdmin 4 pgAdmin 4
pgAdmin 4

2.2 postgresql

```
psql PostgreSQL
Oracle sqlplus
psql
PostgreSQL
```

2.2.1 安装psql数据库

```

psql

```

```
[postgres@pghost1 ~]$ psql postgres postgres
psql (10.0)
Type "help" for help.
postgres=#
```

```
psql postgres postgres
$PGPORT
1921 mydb
pguser mydb tbs_mydb
```

```
- - 创建用户
postgres=# CREATE ROLE pguser WITH ENCRYPTED PASSWORD
'pguser';
```

CREATE ROLE

--创建目录

```
[postgres@pghost1 ~]$ mkdir -p  
/database/pg10/pg_tbs/tbs_mydb
```

--创建表空间

```
postgres=# CREATE TABLESPACE tbs_mydb OWNER pguser LOCATION  
'/database/pg10/pg_tbs/tbs_mydb';  
CREATE TABLESPACE
```

--创建数据库

```
postgres=# CREATE DATABASE mydb  
          WITH OWNER = pguser  
          TEMPLATE = template0  
          ENCODING = 'UTF8'  
          TABLESPACE = tbs_mydb;  
CREATE DATABASE
```

--授权

```
GRANT ALL ON DATABASE mydb TO pguser WITH GRANT OPTION;  
GRANT ALL ON TABLESPACE tbs_mydb TO pguser;
```

```
CREATE DATABASE OWNER pguser  
TEMPLATE template0  
template1 ENCODING UTF8  
TABLESPACE  
psql
```

```
psql [option...] [dbname [username]]
```

```
pgghost1 IP 192.168.28.74  
pgghost2 IP 192.168.28.75 pgghost2
```

pgghost1mydb

```
[postgres@pgghost2 ~]$ psql -h 192.168.28.74 -p 1921 mydb
pguser
Password for user pguser:
psql (10.0)
Type "help" for help.
```

psql\qCTRL+D

```
[postgres@pgghost1 ~]$ psql mydb pguser
psql (10.0)
Type "help" for help.
mydb=> \q
```

psql

2.2.2 psql

psql

```
postgres=# \l
```

List of databases				
Name	Owner	Encoding	Collate	Ctype
Access privileges				

```

-----+-----+-----+-----+-----+-----
mydb      | postgres | UTF8     | C       | C       |
=Tc/postgres      +
|           |           |           |         |         |
postgres=CTc/postgres+
|           |           |           |         |         |
pguser=C/postgres
postgres  | postgres | UTF8     | C       | C       |
template0 | postgres | UTF8     | C       | C       |
=c/postgres      +
|           |           |           |         |         |
postgres=CTc/postgres
template1 | postgres | UTF8     | C       | C       |
=c/postgres      +
|           |           |           |         |         |
postgres=CTc/postgres
(4 rows)

```

1..\db

```

postgres=# \db
              List of tablespaces
   Name   | Owner   | Location
-----+-----+-----
pg_default | postgres |
pg_global  | postgres |
tbs_mydb   | pguser   | /database/pg10/pg_tbs/tbs_mydb
(3 rows)

```

2.\d

```

mydb=> CREATE TABLE test_1(id int4,name text,
        create_time timestamp without time zone default
        clock_timestamp());
CREATE TABLE

mydb=> ALTER TABLE test_1 ADD PRIMARY KEY (id);
ALTER TABLE

```

generate_series(1,50000000) test_1\d

```

mydb=> \d test_1
                                Table "pguser.test_1"
  Column |          Type          | Collation |
Nullable |          Default        |           |
-----+-----+-----+-----+
id       | integer                |           | not
null    |                          |           |
name     | text                   |           |
create_time| timestamp without time zone |           |
clock_timestamp()
Indexes:
    "test_1_pkey" PRIMARY KEY , btree (id)

```

3. 插入数据

向test_1表插入500万条数据

```

mydb=> INSERT INTO test_1(id,name)
        SELECT n,n || '_francs'
        FROM generate_series(1,5000000) n;
INSERT 0 5000000

```

pgbench\dt+pgbenchデータベース

```
mydb=> \dt+ test_1
              List of relations
Schema | Name   | Type  | Owner  | Size  | Description
-----+-----+-----+-----+-----+-----
pguser | test_1 | table | pguser | 287 MB |
(1 row)
```

pgbench\di+pgbenchデータベース

```
mydb=> \di+ test_1_pkey
              List of relations
Schema | Name          | Type  | Owner  | Table  | Size  |
Description
-----+-----+-----+-----+-----+-----+-----
pguser | test_1_pkey   | index | pguser | test_1 | 107 MB |
(1 row)
```

4.\sfpgbenchデータベース

pgbench\sfpbenchデータベース

```
mydb=> \sf random_range
CREATE OR REPLACE FUNCTION pguser.random_range(integer,
integer)
    RETURNS integer
    LANGUAGE sql
AS $function$
    SELECT ($1 + FLOOR(($2 - $1 + 1) * random()))::int4;
$function$
```

`\sf` 显示特殊字符
`random_range(integer, integer)`
PostgreSQL 特殊字符
`\sf` 显示特殊字符

5. \x 显示十六进制

`\x` 显示十六进制

```
mydb=> SELECT * FROM test_1 LIMIT 1;
   id |   name   |          create_time
-----+-----+-----
    1 | 1_pguser | 2017-07-22 11:16:15.97559
(1 row)
```

```
mydb=> \x
Expanded display is on.
mydb=> SELECT * FROM test_1 LIMIT 1;
-[ RECORD 1 ]-----
id           | 1
name         | 1_francs
create_time  | 2017-07-22 11:16:15.97559
```

6. 使用 psql 连接 SQL

`psql` 连接数据库
`psql -E` 显示特殊字符

```
[postgres@pghost1 ~]$ psql -E mydb pguser
psql (10.0)
Type "help" for help.
```

```

mydb=> \db
***** QUERY *****
SELECT spcname AS "Name",
pg_catalog.pg_get_userbyid(spcowner) AS "Owner",
pg_catalog.pg_tablespace_location(oid) AS "Location"
FROM pg_catalog.pg_tablespace
ORDER BY 1;
*****

```

List of tablespaces		
Name	Owner	Location
pg_default	postgres	
pg_global	postgres	
tbs_mydb	pguser	/database/pg10/pg_tbs/tbs_mydb

(3 rows)

7. \???

PostgreSQL \???

mydb=> \?	
General	
\copyright	show PostgreSQL usage and
distribution terms	
\crosstabview [COLUMNS]	execute query and display
results in crosstab	
\errverbose	show most recent error message
at maximum verbosity	
\g [FILE] or ;	execute query (and send results
to file or pipe)	
\gexec	execute query, then execute each
value in its result	
\gset [PREFIX]	execute query and store results
in psql variables	
\gx [FILE]	as \g, but forces expanded

output mode	
\q	quit psql
\watch [SEC]	execute query every SEC seconds
Help	
\? [commands]	show help on backslash commands
\? options	show help on psql command-line
options	
\? variables	show help on special variables
\h [NAME]	help on syntax of SQL commands,
* for all commands	

```

\pset pager on
\pset pager_page 1

```

8. \? HELP

```

psql> \?
\h SQL
\h CREATE TABLESPACE

```

```

postgres=# \h CREATE TABLESPACE
Command:      CREATE TABLESPACE
Description:  define a new tablespace
Syntax:
CREATE TABLESPACE tablespace_name
    [ OWNER { new_owner | CURRENT_USER | SESSION_USER } ]
    LOCATION 'directory'
    [ WITH ( tablespace_option = value [, ... ] ) ]

```

```

\h SQL
SQL

```

2.2.3 psqlコマンド

psqlコマンドを実行すると、標準入力から標準出力へデータをコピーします。
COPYコマンドは、標準入力から標準出力へデータをコピーします。

1. COPYコマンドでSQLコマンドを実行する

2. COPYコマンドでSUPERUSERコマンドを実行する
stdinとstdoutの両方からデータをコピーする
SUPERUSERコマンド

3. COPYコマンドで標準入力から標準出力へデータをコピーする
psqlコマンドを実行する

4. コピーしたデータを標準出力へ書き出す
COPYコマンドを実行する

1. COPYコマンド

COPYコマンドでデータをコピーする
mydbデータベースにtest_copyテーブルを作成する

```
mydb=> CREATE TABLE test_copy(id int4,name text);  
CREATE TABLE
```

test_copy_in.txt TAB
test-copy-in.txt

```
[pg10@pghost1 script]$ cat test_copy_in.txt
1      a
2      b
3      c
```

postgres mydb
test_copy_in.txt test_copy

```
[pg10@pghost1 script]$ psql mydb postgres
psql (10.0)
Type "help" for help.
```

```
mydb=# COPY pguser.test_copy FROM
'/home/postgres/script/test_copy_in.txt';
COPY 3
mydb=# SELECT * FROM pguser.test_copy ;
   id | name
-----+-----
    1 | a
    2 | b
    3 | c
(3 rows)
```

pguser

```
[pg10@pghost1 script]$ psql mydb pguser
psql (10.0)
```

Type "help" for help.

```
mydb=> COPY test_copy FROM
'/home/postgres/script/test_copy_in.txt';
ERROR:  must be superuser to COPY to or from a file
HINT:  Anyone can COPY to stdout or from stdin. psql's \copy
command also works for anyone.
```

```
psql (10.0)
Type "help" for help.
mydb=# COPY test_copy
FROM postgres mydb

```

```
[pg10@pghost1 script]$ psql mydb postgres
psql (10.0)
Type "help" for help.
mydb=# COPY pguser.test_copy TO
'/home/postgres/test_copy.txt';
COPY 3
```

```
test_copy.txt
```

```
[postgres@pghost1 ~]$ cat test_copy.txt
1      a
2      b
3      c
```

```
psql (10.0)
```

```
[postgres@pghost1 ~]$ psql mydb pguser
psql (10.0)
```

Type "help" for help.

```
mydb=> COPY test_copy T0 stdout;
```

1	a
2	b
3	c

[illegible]

DBA
CSV

```
[postgres@pghost1 ~]$ psql mydb postgres
psql (10.0)
Type "help" for help.
```

```
mydb=# COPY pguser.test_copy TO
'/home/postgres/test_copy.csv' WITH csv header;
COPY 4
```

```

with csv header csv
csv office excel
test_copy ID 1

```

```
mydb=# COPY (SELECT * FROM pguser.test_copy WHERE id=1) TO
'/home/postgres/1.txt';
COPY 1
mydb=# \q
[postgres@pghost1 ~]$ cat 1.txt
1      a
```

COPY
<https://www.postgresql.org/docs/10/static/sql-copy.html>

2. \copy

COPY
psql
pgghost2
pguser
pgghost1
mydb
\copy
test_copy

```
[postgres@pgghost2 ~]$ psql -h 192.168.28.74 -p 1921 mydb
pguser
Password for user pguser:
psql (10.0)
Type "help" for help.

mydb=> \copy test_copy to '/home/postgres/test_copy.txt';
COPY 3
```

test_copy.txt

```
[postgres@pgghost2 ~]$ cat test_copy.txt
1      a
2      b
3      c
```

```
\copy test_copy_in.txt
```

```
[postgres@pghost2 ~]$ cat test_copy_in.txt
4      d
```

```
\copy test_copy_in.txt
```

```
[postgres@pghost2 ~]$ psql -h 192.168.28.74 -p 1921 mydb
pguser
Password for user pguser:
psql (10.0)
Type "help" for help.
```

```
mydb=> \copy test_copy from
'/home/postgres/test_copy_in.txt';
COPY 1
mydb=> SELECT * FROM test_copy WHERE id=4;
   id | name
-----+-----
    4 | d
(1 row)
```

```
\copy
COPY
```

2.2.4 psql

psql

```
psql [option...] [dbname [username]]
```

dbname username
option

1.-A

psql SQL

```
[postgres@pghost1 ~]$ psql -c "SELECT * FROM user_ini WHERE
id=1" mydb pguser
   id | user_id | user_name |          create_time
-----+-----+-----+-----
    1 | 186536 | KTU89H    | 2017-08-05 15:59:25.359148+08
(1 row)
--
```

psql -A

```
[postgres@pghost1 ~]$ psql -A -c "SELECT * FROM user_ini
WHERE id=1" mydb pguser
id|user_id|user_name|create_time
1|186536|KTU89H|2017-08-05 15:59:25.359148+08
(1 row)
--
```

psql -A 显示所有列名
psql -t 显示一行

2. -t 显示一行

psql 显示一行 -t -t 显示一行
psql 显示一行

```
[postgres@pghost1 ~]$ psql -t -c "SELECT * FROM user_ini  
WHERE id=1" mydb pguser  
1 | 186536 | KTU89H | 2017-08-05 15:59:25.359148+08  
--
```

psql 显示一行 -t -A 显示所有列名
psql 显示一行

```
[postgres@pghost1 ~]$ psql -At -c "SELECT * FROM user_ini  
WHERE id=1" mydb pguser  
1|186536|KTU89H|2017-08-05 15:59:25.359148+08
```

psql 显示一行 shell 显示一行
psql 显示一行

```
[postgres@pghost1 ~]$ psql -At -c "SELECT user_name FROM  
user_ini WHERE id=1" mydb pguser  
KTU89H
```

3.-q

psql SQL -q
test_q.sql SQL

```
DROP TABLE if exists test_q;  
CREATE TABLE test_q(id int4);  
TRUNCATE TABLE test_q;  
INSERT INTO test_q values (1);  
INSERT INTO test_q values (2);
```

test_q.sql

```
[postgres@pghost1 ~]$ psql mydb pguser -f test_q.sql  
DROP TABLE  
CREATE TABLE  
TRUNCATE TABLE  
INSERT 0 1  
INSERT 0 1
```

test_q.sql -q

```
[postgres@pghost1 ~]$ psql -q mydb pguser -f test_q.sql  
--
```

-q -c -f

2.2.5 psqlとsql

psql -c 実行するSQLコマンドを指定して実行する

```
[postgres@pghost1 ~]$ psql -c "SELECT current_user;"
current_user
-----
postgres
(1 row)
```

-c 実行するSQLコマンドを指定して実行する
psql -At 実行するSQLコマンドを指定して実行する

```
[postgres@pghost1 ~]$ psql -At -c "SELECT current_user;"
postgres
```

psql SQLコマンドを実行する
test_2.sql

```
CREATE TABLE test_2(id int4);
INSERT INTO test_2 VALUES (1);
INSERT INTO test_2 VALUES (2);
INSERT INTO test_2 VALUES (3);
```

-f 実行するSQLコマンドを指定して実行する

```
[postgres@pgghost1 ~]$ psql mydb pguser -f script/test_2.sql
CREATE TABLE
INSERT 0 1
INSERT 0 1
INSERT 0 1
```

データベースの操作はSQLで実行する



psql-single-transaction-1
データベースの操作はSQLで実行する
SQLでデータベースの操作を行う

2.2.6 psqlでデータベースのSQL

psqlでデータベースのSQLを実行する

```
SELECT * FROM table_name WHERE column_name = 1
```

データベースの操作はSQLで実行する

1..\setでデータベースの操作

\setでデータベースの操作を行う
valueでデータベースの操作を行う

```
\set name value
```

```
test_copy v_id 2 id
2
```

```
mydb=> \set v_id 2
mydb=> SELECT * FROM test_copy WHERE id=:v_id;
   id | name
-----+-----
    2 | b
(1 row)
```

```
\set

```

```
mydb=> \set v_id
```

```
\set pgbench
\set
```

2.psql-v

```
psql-v
select_1.sql
```

```
SELECT * FROM test_3 WHERE id=:v_id;
```

psql -v select_1.sql

```
[postgres@pghost1 ~]$ psql -v v_id=1 mydb pguser -f
select_1.sql
   id | name
-----+-----
     1 | a
(1 row)
```

v_id=1

2.2.7 psql

DBA

psql

1.

.psqlrc psql -X psql

~/.psqlrc

SQL

```
SELECT pid,username,datname,query,client_addr
FROM pg_stat_activity
```

```
WHERE pid <> pg_backend_pid() AND state='active' ORDER BY
query;
```

pg_stat_activity PostgreSQL 表
pid username
datname query
SQL state active query
SQL client_addr IP state
表

·active PostgreSQL

·idle

·idle in transaction
SQL

·idle in transaction aborted idle in
transaction SQL

<https://www.postgresql.org/docs/10/static/monitoring-stats.html#pg-stat-activity-view>

~/.psqlrc
\set SQL

```
--\set active_session 'select
pid,username,datname,query,client_addr from pg_stat_activity
where pid <> pg_backend_pid() and state=\'active\' order by
query;'
```

active_session

```
postgres=# :active_session
      pid | username | datname | query
-----+-----+-----+-----
      14351 | pguser  | mydb    | update test_per1 set
create_time=now() WHERE id=$1; |
      14352 | pguser  | mydb    | update test_per1 set
create_time=now() WHERE id=$1; |
      14353 | pguser  | mydb    | update test_per1 set
create_time=now() WHERE id=$1; |
      14354 | pguser  | mydb    | update test_per1 set
create_time=now() WHERE id=$1; |
      14355 | pguser  | mydb    | update test_per1 set
create_time=now() WHERE id=$1; |
(5 rows)
```

SQLactive_session

2.

PostgreSQL
SQL

```
SELECT
pid,username,datname,query,client_addr,wait_event_type,wait_e
vent
FROM pg_stat_activity
WHERE pid <> pg_backend_pid() AND wait_event is not null
ORDER BY wait_event_type;
```

```
~~~~~\set~~~~~ ~/.psqlrc~~~~~
\set~~~SQL~~~~~
```

```
--~~~~~
\set wait_event 'select
pid,username,datname,query,client_addr,wait_event_type,wait_e
vent from pg_stat_activity where pid <> pg_backend_pid() and
wait_event is not null order by wait_event_type;
```

```
~~~~~wait_event~~~~~
~~~~~
```

```
postgres=# :wait_event
      pid | username | datname | query | client_addr |
wait_event_type |      wait_event
-----+-----+-----+-----+-----+-----
-----+-----
      2652 |          |         |       |              | Activity
| AutoVacuumMain
      2655 | postgres |         |       |              |
Activity | LogicalLauncherMain
      2650 |          |         |       |              |
Activity | BgWriter Hibernate
      2649 |          |         |       |              |
Activity | CheckpointerMain
      2651 |          |         |       |              | Activity
| WalWriterMain
(5 rows)
```



```
count
-----
1000000
(1 row)

Time: 47.114 ms
```

count 47.114ms
SQL \timing
timing

```
mydb=> \timing
Timing is off.
```

2. \watch SQL

\watch SQL
SQL

```
\watch [ seconds ]
```

seconds 2
now

```
mydb=> SELECT now();
now
-----
2017-08-14 11:20:02.157567+08
(1 row)
```



```
mydb=> ALTER TABLE test_1 DROP C0
COLUMN          CONSTRAINT
mydb=> ALTER TABLE test_1 DROP C0
```

4. 使用 psql 连接 PostgreSQL

psql 是 PostgreSQL 的交互式 SQL 客户端。

```
[postgres@pghost1 ~]$ psql mydb pguser
psql (10.0)
Type "help" for help.
```

```
mydb=> SELECT count(*) FROM pg_stat_activity ; -- 显示正在运行的进程
      count
      ----
         1
```

使用 psql 连接 PostgreSQL 数据库时，可以通过以下命令指定连接参数：

PostgreSQL 支持 readline 库，用于提供行编辑功能。

PostgreSQL 还支持 `--without-readline` 选项，用于禁用 readline 库。

5. psql 连接 PostgreSQL

使用 psql 连接 PostgreSQL 数据库时，可以通过以下命令指定连接参数：

```
psql postgres=#
```

```
[postgres@pghost1 ~]$ psql
psql (10.0)
Type "help" for help.
```

```
postgres=#
```

psql psql

.%M psql -h
Unix [local]

.%>

.%n
SET SESSION AUTHORIZATION

.%/

.%# " #" ">"
SET SESSION
AUTHORIZATION

.%p

.%R PROMPT1 "="
"

psql prompt1

```
[postgres@pghost1 ~]$ psql
psql (10.0)
Type "help" for help.
```

```
postgres=# \echo :PROMPT1
%/%R%#
```

```
    \echo:PROMPT1
PROMPT1psql
%/%R%#
postgres%R"=" %#"#
%MUnix
PROMPT1%M%R%#
```

```
[postgres@pghost1 ~]$ psql
psql (10.0)
Type "help" for help.
```

```
postgres=# \set PROMPT1 '%M%R%#'
[local]=#
```

```
psql[local]pghost2
pghost1PROMPT1
"%M%R%#"
```

```
[postgres@pghost2 ~]$ psql -h 192.168.28.74 mydb pguser -p
1921
Password for user pguser:
psql (10.0)
Type "help" for help.
```

```
mydb=> \set PROMPT1 '%M%R%#'
192.168.28.74=>
```

PROMPT1 "%M%R%#" 192.168.28.74 IP psql -h

```
[postgres@pghost2 ~]$ psql -h 192.168.28.74 mydb pguser -p 1921
Password for user pguser:
psql (10.0)
Type "help" for help.

mydb=> \set PROMPT1 '%/@%M:%>%R%#'
mydb@192.168.28.74:1921=>
```

PROMPT1 "%/@%M %>%R%#" "%>" PROMPT1 .psqlrc .psqlrc 2.2.7 .psqlrc

```
[postgres@pghost2 ~]$ touch .psqlrc
[postgres@pghost2 ~]$ vim .psqlrc
\set PROMPT1 '%/@%M:%>%R%#'
```

```
[postgres@pghost2 ~]$ psql -h 192.168.28.74 mydb pguser -p 1921
Password for user pguser:
psql (10.0)
```


Type "help" for help.

mydb@192.168.28.74:1921=>

PROMPT1 .psqlrc
psql .psqlrc



psql PROMPT1
PROMPT2 PROMPT3 PROMPT1 psql
PROMPT2
PROMPT2
PROMPT2
PROMPT1 PROMPT3
SQLCOPYFROMSTDIN

2.3 四角

```

root@kali:~# sudo su postgres
postgres@kali:~$ psql
psql (15.3)
Type "help" for help.

postgres=# \du
               List of databases
   Name | Owner | Encoding | Collate | Ctype | Tablespace |
-----+-----+-----+-----+-----+-----+
 postgres | postgres | UTF8     | en_US.UTF-8 | en_US.UTF-8 |          |
 template0 | postgres | UTF8     | en_US.UTF-8 | en_US.UTF-8 |          |
 template1 | postgres | UTF8     | en_US.UTF-8 | en_US.UTF-8 |          |
(3 rows)

postgres=# \c postgres
You are now connected to database "postgres" as user "postgres".
postgres=#

```

3

```

    postgresql postgresql
    postgresql postgresql postgresql postgresql postgresql postgresql postgresql postgresql
    postgresql postgresql postgresql postgresql/postgresql postgresql postgresql postgresql
    postgresql postgresql postgresql postgresql json/jsonb postgresql postgresql
    postgresql postgresql postgresql postgresql postgresql postgresql postgresql postgresql

```

3.1 数据类型

PostgreSQL数据类型分为基本数据类型和复合数据类型。基本数据类型包括integer、smallint、serial等。

3.1.1 基本数据类型

PostgreSQL基本数据类型3-1

表 3-1 基本数据类型

类型名称	存储长度	描述	范围
smallint	2 字节	小范围整数类型	32 768 到 +32 767
integer	4 字节	整数类型	-2 147 483 648 到 +2 147 483 647
bigint	8 字节	大范围整数类型	-9 223 372 036 854 775 808 到 +9 223 372 036 854 775 807
decimal	可变	用户指定精度	小数点前 131 072 位；小数点后 16 383 位
numeric	可变	用户指定精度	小数点前 131 072 位；小数点后 16 383 位
real	4 字节	变长，不精确	6 位十进制精度

(续)

类型名称	存储长度	描述	范围
double precision	8 字节	变长，不精确	15 位十进制精度
smallserial	2 字节	smallint 自增序列	1 到 32 767
serial	4 字节	integer 自增序列	1 到 2 147 483 647
bigserial	8 字节	bigint 自增序列	1 到 9 223 372 036 854 775 807

smallint、integer、bigint 都是整数类型。smallint 是 2 字节整数，integer 是 4 字节整数，bigint 是 8 字节整数。

integer int8
integer integer
bigint integer

```
mydb=> CREATE TABLE test_integer (id1 integer,id2 int4) ;  
CREATE TABLE
```

decimal numeric
numeric numeric
numeric

```
NUMERIC(precision, scale)
```

precision numeric scale
18.222 precision 5 scale
3 precision scale 0
numeric
numeric

real double precision real
4 double precision 8
real

smallserial serial bigserial
serial
test_serial id serial

```
mydb=> CREATE TABLE test_serial (id serial,flag text);
CREATE TABLE
```

test_serial

```
mydb=> \d test_serial
```

Column	Type	Collation	Nullable
id	integer		not null
flag	text		

Table "pguser.test_serial"

nextval('test_serial_id_seq'::regclass)

test_serial_id_seq

serial

```
mydb=> INSERT INTO test_serial(flag) VALUES ('a');
INSERT 0 1
mydb=> INSERT INTO test_serial(flag) VALUES ('b');
INSERT 0 1
mydb=> INSERT INTO test_serial(flag) VALUES ('c');
INSERT 0 1
mydb=> SELECT * FROM test_serial;
```

id	flag
1	a
2	b
3	c

(3 rows)

3.1.2

PostgreSQL 算術演算関数

四則演算

```
mydb=> SELECT 1+2,2*3,4/2,8%3;
      ?column? | ?column? | ?column? | ?column?
-----+-----+-----+-----
              3 |          6 |          2 |          2
```

剰余演算

```
mydb=> SELECT mod(8,3);
      mod
-----
      2
(1 row)
```

四捨五入

```
mydb=> SELECT round(10.2),round(10.9);
      round | round
-----+-----
          10 |          11
(1 row)
```

上取捨

```
mydb=> SELECT ceil(3.6),ceil(-3.6);
      ceil | ceil
-----+-----
         4 |        -3
(1 row)
```

□□□□□□□□□□□□□□□□□□□□□□

```
mydb=> SELECT floor(3.6), floor(-3.6);
      floor | floor
-----+-----
          3 |      -4
(1 row)
```

3.2 字符串

本章介绍PostgreSQL字符串数据类型及其使用。

3.2.1 字符串类型

PostgreSQL字符串类型如表3-2所示。

表3-2 字符串类型

字符类型名称	描 述
character varying(n), varchar(n)	变长，字符最大数有限制
character(n), char(n)	定长，字符数没达到最大值则使用空白填充
text	变长，无长度限制

character varying(n) 变长字符串，n 为最大长度。
character(n) 定长字符串，n 为长度，不足 n 个字符时用空格填充。
text 变长字符串，无长度限制。

```
mydb=> CREATE TABLE test_char(col1 varchar (4),col2
character(4));
CREATE TABLE
mydb=> INSERT INTO test_char(col1,col2) VALUES ('a','a');
INSERT 0 1
```

test_char col1 character varying
4 col2 character 4

```
mydb=> SELECT char_length(col1),char_length(col2) FROM
test_char ;
  char_length | char_length
-----+-----
          1 |           1
(1 row)
```

char_length string
1

```
mydb=> SELECT octet_length(col1),octet_length(col2) FROM
test_char ;
  octet_length | octet_length
-----+-----
          1 |           4
(1 row)
```

octet_length string col2
4 col2 character

character varying n
character n
character 1

位置取得関数 position

```
mydb=> SELECT position('a' in 'abcd');
position
-----
          1
(1 row)
```

文字列抽出関数 substring

```
mydb=> SELECT substring('francs' from 3 for 4);
substring
-----
      ancs
(1 row)
```

文字列分割関数 split_part

split_part(string text, delimiter text, field int)

delimiterは文字列、stringは分割対象の文字列、fieldは分割するフィールド番号(1から開始)

```
mydb=> SELECT split_part('abc@def1@nb','@',2);
split_part
-----
      def1
(1 row)
```

3.3 日期/时间

PostgreSQL日期/时间数据类型
PostgreSQL日期/时间数据类型

3.3.1 日期/时间

PostgreSQL日期/时间数据类型3-3

表3-3 日期/时间数据类型

字符类型名称	存储长度	描述
timestamp [(p)] [without time zone]	8 字节	包括日期和时间，不带时区，简写成 timestamp
timestamp [(p)] with time zone	8 字节	包括日期和时间，带时区，简写成 timestamptz
date	4 字节	日期，但不包含一天中的时间
time [(p)] [without time zone]	8 字节	一天中的时间，不包含日期，不带时区
time [(p)] with time zone	12 字节	一天中的时间，不包含日期，带时区
interval [fields] [(p)]	16 字节	时间间隔

PostgreSQL日期/时间数据类型
now()now()now()
timestamp(p)with time zone

```
mydb=> SELECT now();
          now
```

```
-----
2017-07-29 09:44:25.493425+08
(1 row)
```

timestamp without time zone

```
mydb=> SELECT now()::timestamp without time zone;
              now
-----
2017-07-29 09:44:55.804403
(1 row)
```

date

```
mydb=> SELECT now()::date;
              now
-----
2017-07-29
(1 row)
```

time without time zone

```
mydb=> SELECT now()::time without time zone;
              now
-----
09:45:49.390428
(1 row)
```

time with time zone

```
mydb=> SELECT now()::
time with time zone;
```

```

              now
-----
09:45:57.13139+08
(1 row)

```

interval hour day
month year

```

mydb=> SELECT now(),now()+interval'1 day';
              now                |                ?column?
-----+-----
2017-07-29 09:47:26.026418+08 | 2017-07-30
09:47:26.026418+08
(1 row)

```

timestamp(0)
precision
6
0

```

mydb=> SELECT now(), now()::timestamp(0);
              now                |                now
-----+-----
2017-07-29 09:59:42.688445+08 | 2017-07-29 09:59:43
(1 row)

```

3.3.2 日期/时间戳

timestamp(0)

□□□□□□□□□□

```
mydb=> SELECT date '2017-07-29' + interval'1 days';
        ?column?
-----
2017-07-30 00:00:00
(1 row)
```

□□□□□□□□□□

```
mydb=> SELECT date '2017-07-29' - interval'1 hour';
        ?column?
-----
2017-07-28 23:00:00
(1 row)
```

□□□□□□□□□□

```
mydb=> SELECT 100* interval '1 second';
        ?column?
-----
00:01:40
(1 row)
```

□□□□□□□□□□

```
mydb=> SELECT interval '1 hour' / double precision '3';
        ?column?
-----
00:20:00
(1 row)
```

timestamp

```
mydb=> SELECT EXTRACT( month FROM now()),EXTRACT(day FROM now());
```

date_part		date_part
-----	+	-----
7		29

(1 row)

```
mydb=> SELECT EXTRACT( hour FROM now()), extract (minute FROM now());
```

date_part		date_part
-----	+	-----
11		14

```
mydb=> SELECT EXTRACT( second FROM now());
```

date_part

43.031366

(1 row)

```
mydb=> SELECT EXTRACT( week FROM now());
```

date_part

30

(1 row)

□□□□□□□□□□□□□□□□

```
mydb=> SELECT EXTRACT( doy FROM now());
      date_part
-----
          210
(1 row)
```

3.4 数据类型

在 PostgreSQL 中，数据类型分为基本数据类型和复合数据类型。基本数据类型包括整数、浮点数、字符串、日期、时间、布尔值等。复合数据类型包括数组、JSON、JSONB 等。在本节中，我们将重点介绍 PostgreSQL 中的布尔数据类型。图 3-4 展示了布尔数据类型的存储长度和描述。

图 3-4 布尔数据类型

字符类型名称	存储长度	描述
boolean	1 字节	状态为 true 或 false

在 PostgreSQL 中，布尔数据类型可以使用以下值进行赋值：
true 或 TRUE 或 t 或 true 或 y 或 yes 或 on 或 1
false 或 FALSE 或 f 或 false 或 n 或 no 或 off 或 0

```
mydb=> CREATE TABLE test_boolean(cola boolean,colb boolean);
CREATE TABLE
mydb=> INSERT INTO test_boolean (cola,colb) VALUES
('true','false');
INSERT 0 1
mydb=> INSERT INTO test_boolean (cola,colb) VALUES ('t','f');
INSERT 0 1
mydb=> INSERT INTO test_boolean (cola,colb) VALUES
('TRUE','FALSE');
INSERT 0 1
mydb=> INSERT INTO test_boolean (cola,colb) VALUES
('yes','no');
INSERT 0 1
mydb=> INSERT INTO test_boolean (cola,colb) VALUES ('y','n');
INSERT 0 1
```


3.5 数据类型

PostgreSQL 支持多种数据类型，包括 IP 地址、网络地址、MAC 地址等。PostgreSQL 支持 IPv4、IPv6、MAC 地址等数据类型。PostgreSQL 支持 IP 地址、网络地址、MAC 地址等数据类型。

3.5.1 网络地址

PostgreSQL 支持多种网络地址数据类型，包括 IPv4、IPv6 等。

表 3-5 网络地址数据类型

字符类型名称	存储长度	描述
cidr	7 或 19 字节	IPv4 和 IPv6 网络
inet	7 或 19 字节	IPv4 和 IPv6 网络
macaddr	6 字节	MAC 地址
macaddr8	8 字节	MAC 地址（EUI-64 格式）

inet、cidr 数据类型用于存储 IP 地址和网络地址。address/y 数据类型用于存储 IP 地址和网络地址。address 数据类型用于存储 IP 地址，支持 IPv4 和 IPv6。address/y 数据类型用于存储 IP 地址和网络地址，支持 IPv4 和 IPv6。address/y 数据类型用于存储 IP 地址和网络地址，支持 IPv4 和 IPv6。

inet、cidr 数据类型用于存储 IP 地址和网络地址。address/y 数据类型用于存储 IP 地址和网络地址。address 数据类型用于存储 IP 地址，支持 IPv4 和 IPv6。address/y 数据类型用于存储 IP 地址和网络地址，支持 IPv4 和 IPv6。

```
mydb=> SELECT '192.168.2.1000'::inet;
ERROR:  invalid input syntax for type inet: "192.168.2.1000"
LINE 1: select '192.168.2.1000'::inet;
```

inet cidr

1 cidr inet

```
mydb=> SELECT '192.168.1.100'::cidr;
      cidr
-----
 192.168.1.100/32
(1 row)

mydb=> SELECT '192.168.1.100/32'::inet;
      inet
-----
 192.168.1.100
(1 row)

mydb=> SELECT '192.168.0.0/16'::inet;
      inet
-----
 192.168.0.0/16
(1 row)
```

2 cidr IP inet

```
mydb=> SELECT '192.168.2.0/8'::cidr;
ERROR:  invalid cidr value: "192.168.2.0/8"
LINE 1: select '192.168.2.0/8'::cidr;
DETAIL:  Value has bits set to right of mask.
```

```
mydb=> SELECT '192.168.2.0/8'::inet;  
      inet
```

```
-----  
      192.168.2.0/8  
(1 row)
```

```
mydb=> SELECT '192.168.2.0/24'::cidr;  
      cidr
```

```
-----  
      192.168.2.0/24  
(1 row)
```

cidr inet macaddr
macaddr8 MAC

3.5.2

PostgreSQL 3-6

3-6

操 作 符	描 述	举 例
<	小于	inet '192.168.1.5' < inet '192.168.1.6'
<=	小于等于	inet '192.168.1.5' <= inet '192.168.1.5'
=	等于	inet '192.168.1.5' = inet '192.168.1.5'
>=	大于等于	inet '192.168.1.5' >= inet '192.168.1.5'
>	大于	inet '192.168.1.5' > inet '192.168.1.4'
<>	不等于	inet '192.168.1.5' <> inet '192.168.1.4'
<<	被包含	inet '192.168.1.5' << inet '192.168.1/24'
<<=	被包含或等于	inet '192.168.1/24' <<= inet '192.168.1/24'
>>	包含	inet '192.168.1/24' >> inet '192.168.1.5'
>>=	包含或等于	inet '192.168.1/24' >>= inet '192.168.1/24'
&&	包含或被包含	inet '192.168.1/24' && inet '192.168.1.80/28'
~	按位取反	~ inet '192.168.1.6'
&	按位与	inet '192.168.1.6' & inet '0.0.0.255'
	按位或	inet '192.168.1.6' inet '0.0.0.255'
+	加	inet '192.168.1.6' + 25
-	减	inet '192.168.1.43' - 36
-	减	inet '192.168.1.43' - inet '192.168.1.19'

3.5.3 子网掩码

PostgreSQL子网掩码

IP子网掩码

```
mydb=> SELECT host(cidr '192.168.1.0/24');
          host
-----
192.168.1.0
(1 row)
```

IP子网掩码

```
mydb=> SELECT text(cidr '192.168.1.0/24');
      text
```

```
-----
 192.168.1.0/24
(1 row)
```

```
□□□□□□□□□□□□□□□□□□□□□□□□
```

```
mydb=> SELECT netmask(cidr '192.168.1.0/24');
      netmask
```

```
-----
255.255.255.0
```

3.6 数组

PostgreSQL 支持一维数组，支持多种数据类型的一维数组，包括整数、文本、浮点数、日期、时间、枚举、JSON 等。

3.6.1 创建数组表

创建数组表，使用 `CREATE TABLE` 语句，指定表名、列名和数据类型。数组列的数据类型必须在 `CREATE TABLE` 语句中指定，且必须使用 `array` 关键字。

```
CREATE TABLE test_array1 (  
    id          integer,  
    array_i     integer[],  
    array_t     text[]  
);
```

在 `integer[]` 和 `integer` 之间使用 `array` 关键字，在 `text[]` 和 `text` 之间使用 `array` 关键字。

3.6.2 插入数组数据

插入数组数据，使用 `INSERT` 语句，指定表名、列名和数据值。数组列的数据值必须使用花括号括起来，并且使用分隔符分隔。

```
{ val1 delim val2 delim ... }
```

CREATE TABLE test_array1 (id INT, array_i JSON, array_t JSON, array ARRAY) DELIMITED ROWS BY VALUES;

```
mydb=> SELECT '{1,2,3}';
      ?column?
-----
      {1,2,3}
(1 row)
```

test_array1

```
mydb=> INSERT INTO test_array1(id,array_i,array_t)
VALUES (1,'{1,2,3}','{"a","b","c"}');
INSERT 0 1
```

test_array1 ARRAY

```
mydb=> SELECT array[1,2,3];
      array
-----
      {1,2,3}
(1 row)
```

test_array2

```
mydb=> INSERT INTO test_array1(id,array_i,array_t)
VALUES (2,array[4,5,6],array['d','e','f']);
INSERT 0 1
```

test_array2

```
mydb=> SELECT * FROM test_array1;
  id | array_i | array_t
-----+-----+-----
    1 | {1,2,3} | {a,b,c}
    2 | {4,5,6} | {d,e,f}
(2 rows)
```

3.6.3

```
mydb=> SELECT array_i FROM test_array1 WHERE id=1;
 array_i
-----
 {1,2,3}
(1 row)
```

```
mydb=> SELECT array_i[1],array_t[3] FROM test_array1 WHERE
id=1;
 array_i | array_t
-----+-----
        1 | c
(1 row)
```

3.6.4

PostgreSQL array_append

```
array_append(anyarray, anyelement)
```

array_append

```
mydb=> SELECT array_append(array[1,2,3],4);
       array_append
-----
      {1,2,3,4}
(1 row)
```

array_append ||

```
mydb=> SELECT array[1,2,3] || 4;
       ?column?
-----
      {1,2,3,4}
(1 row)
```

array_remove

```
array_remove(anyarray, anyelement)
```

array_remove

```
mydb=> SELECT array[1,2,2,3],array_remove(array[1,2,2,3],2);
       array      | array_remove
-----+-----
 {1,2,2,3} | {1,3}
(1 row)
```

```
mydb=> UPDATE test_array1 SET array_i[3]=4 WHERE id=1 ;
UPDATE 1
```

```
mydb=> UPDATE test_array1 SET array_i=array[7,8,9] WHERE
id=1;
UPDATE 1
```

3.6.5

PostgreSQL

3-7

操 作 符	描 述	举 例	结 果
=	等于	ARRAY[1,1,2,1,3,1]::int[] = ARRAY[1,2,3]	t
<>	不等于	ARRAY[1,2,3] <> ARRAY[1,2,4]	t
<	小于	ARRAY[1,2,3] < ARRAY[1,2,4]	t
>	大于	ARRAY[1,4,3] > ARRAY[1,2,4]	t
<=	小于等于	ARRAY[1,2,3] <= ARRAY[1,2,3]	t
>=	大于等于	ARRAY[1,4,3] >= ARRAY[1,4,3]	t
@>	包含	ARRAY[1,4,3] @> ARRAY[3,1]	t
<@	被包含	ARRAY[2,7] <@ ARRAY[1,7,4,2,6]	t
&&	重叠（具有公共元素）	ARRAY[1,4,3] && ARRAY[2,1]	t
	数组和数组串接	ARRAY[1,2,3] ARRAY[4,5,6]	{1,2,3,4,5,6}
	数组和数组串接	ARRAY[1,2,3] ARRAY[[4,5,6],[7,8,9]]	{{1,2,3},{4,5,6},{7,8,9}}
	元素和数组串接	3 ARRAY[4,5,6]	{3,4,5,6}
	数组和元素串接	ARRAY[4,5,6] 7	{4,5,6,7}

3.6.6 数组

PostgreSQL 数组操作函数

```
mydb=> SELECT
array_append(array[1,2],3),array_remove(array[1,2],2);
  array_append | array_remove
-----+-----
    {1,2,3}   | {1}
(1 row)
```

数组维度

```
mydb=> SELECT array_ndims(array[1,2]);
  array_ndims
-----
          1
```


1
(1 row)

array_length

```
mydb=> SELECT array_length(array[1,2],1);
array_length
```

2
(1 row)

array_position

```
mydb=> SELECT array_position(array['a','b','c','d'],'d');
array_position
```

4
(1 row)

array_replace

array_replace(anyarray, anyelement, anyelement)

array_replace(anyarray, anyelement, anyelement)

```
mydb=> SELECT array_replace(array[1,2,5,4],5,10);
array_replace
```

```
{1,2,10,4}
(1 row)
```

array_to_string

array_to_string(anyarray, text [, text])

text text text
NULL

```
mydb=> SELECT array_to_string(array[1,2,null,3],',','10');
         array_to_string
-----
        1,2,10,3
(1 row)
```

3.7 範囲型

範囲型は、特定の範囲の値を格納するためのデータ型です。PostgreSQLでは、整数、数値、時刻、日付などの範囲型が提供されています。

3.7.1 範囲型の定義

PostgreSQLでは、範囲型を定義するための構文が提供されています。

- ・int4range—integer範囲型
- ・int8range—bigint範囲型
- ・numrange—numeric範囲型
- ・tsrange—timestamp範囲型
- ・tstzrange—timestamp with time zone範囲型
- ・daterange—date範囲型

範囲型を定義するSQL文の例として、integer範囲型を定義します。

```
mydb=> SELECT int4range(4,7) @> 4;
      ?column?
-----
      t
(1 row)
```

```
mydb=> SELECT int4range(4,7)@>int4range(4,6);
      ?column?
-----
      t
(1 row)
```

```
mydb=> SELECT int4range(4,7)=int4range(4,6,'[]');
         ?column?
-----
         t
(1 row)
```

3.7.4 五五五五五

lower(int4range(1,10))

```
mydb=> SELECT lower(int4range(1,10));
        lower
-----
         1
(1 row)
```

upper(int4range(1,10))

```
mydb=> SELECT upper(int4range(1,10));
        upper
-----
        10
(1 row)
```

isempty(int4range(1,10))

```
mydb=> SELECT isempty(int4range(1,10));
        isempty
-----
             f
(1 row)
```

3.7.5 GiST indexes

GiST indexes are used for storing and indexing data that can be represented by a tree structure. The data is stored in a B-tree structure, and the index is used to find the data quickly. The data is stored in a B-tree structure, and the index is used to find the data quickly. The data is stored in a B-tree structure, and the index is used to find the data quickly.

```
CREATE INDEX idx_ip_address_range ON ip_address USING gist (  
ip_range);
```

3.8 json/jsonb

PostgreSQL implements the JSON and JSONB data types. JSON is JavaScript Object Notation (JSON) and JSONB is JSON with binary storage. JSON is a text-based format, while JSONB is a binary format. Both are used for storing and querying structured data.

3.8.1 json

PostgreSQL 9.2 introduced the json data type. It is a text-based format for storing and querying structured data. The json data type is used for storing and querying JSON documents. The json data type is implemented as a C type, which allows for efficient storage and retrieval of JSON data.

```
mydb=> SELECT '{"a":1,"b":2}'::json;
      json
-----
{"a":1,"b":2}
```

Creating a table with a json column:

```
mydb=> CREATE TABLE test_json1 (id serial primary key,name
json);
CREATE TABLE
```

Inserting data into the table:

```
mydb=> INSERT INTO test_json1 (name)
VALUES ('{"col1":1,"col2":"francs","col3":"male"}');
INSERT 0 1
```

```
mydb=> INSERT INTO test_json1 (name)
VALUES ('{"col1":2,"col2":"fp","col3":"female"}');
INSERT 0 1
```

test_json1

```
mydb=> SELECT * FROM test_json1;
   id |                               name
-----+-----
    1 | {"col1":1,"col2":"francs","col3":"male"}
    2 | {"col1":2,"col2":"fp","col3":"female"}
```

3.8.2 json

“->” json

```
mydb=> SELECT  name -> 'col2' FROM test_json1 WHERE id=1;
?column?
-----
"francs"
(1 row)
```

json“->>”

```

mydb=> SELECT '{"bar": "baz", "balance": 7.77,
"active":false}':::json;
          json
-----
{"bar": "baz", "balance": 7.77, "active":false}
(1 row)

```

jsonb

```

mydb=> SELECT ' {"id":1,      "name":"francs"}':::jsonb;
          jsonb
-----
{"id": 1, "name": "francs"}
(1 row)

```

id name
json

```

mydb=> SELECT ' {"id":1,      "name":"francs"}':::json;
          json
-----
{"id":1,      "name":"francs"}
(1 row)

```

jsonb

```

mydb=> SELECT ' {"id":1,
"name":"francs",
"remark":"a good guy!",
"name":"test"
}':::jsonb;
          jsonb
-----

```

```
      {"id": 1, "name": "test", "remark": "a good guy!"}
(1 row)
```

name name json

jsonb json

3.8.4 jsonb json

json "->>"

```
mydb=> SELECT  name ->> 'col2' FROM test_json1 WHERE id=1;
      ?column?
-----
      francs
(1 row)
```

```
mydb=> SELECT  '{"a":1, "b":2}'::jsonb ? 'a';
      ?column?
-----
      t
(1 row)
```

json数据类型/操作函数

```
mydb=> SELECT '{"a":1, "b":2}':::jsonb - 'a';
      ?column?
-----
 {"b": 2}
(1 row)
```

3.8.5 jsonb/json

json/jsonb数据类型互相转换

json数据类型/操作函数

```
mydb=> SELECT * FROM json_each('{"a":"foo", "b":"bar"}');
      key | value
-----+-----
      a  | "foo"
      b  | "bar"
(2 rows)
```

json数据类型/操作函数

```
mydb=> SELECT * FROM json_each_text('{"a":"foo",
      "b":"bar"}');
      key | value
-----+-----
      a  | foo
      b  | bar
(2 rows)
```

row_to_json
json
json

```
mydb=> SELECT * FROM test_copy WHERE id=1;
```

```
  id | name  
-----+-----  
   1 | a  
(1 row)
```

```
mydb=> SELECT row_to_json(test_copy) FROM test_copy WHERE  
id=1;
```

```
 row_to_json  
-----  
 {"id":1,"name":"a"}  
(1 row)
```

json

```
mydb=> SELECT * FROM json_object_keys('{"a":"foo",  
"b":"bar"}');
```

```
 json_object_keys  
-----  
 a  
 b  
(2 rows)
```

3.8.6 jsonb/

jsonb/sex/

```
mydb=> SELECT '{"name":"francs","age":"31"}'::jsonb ||  
'{"sex":"male"}'::jsonb;  
?column?
```

```
-----  
{"age": "31", "sex": "male", "name": "francs"}  
(1 row)
```

jsonb/#####“-”#####
#####“#-”#####/#####“-”###/#####
##

```
mydb=> SELECT '{"name": "James", "email":  
"james@localhost"}'::jsonb  
- 'email';  
?column?
```

```
-----  
{"name": "James"}  
(1 row)
```

```
mydb=> SELECT '["red","green","blue"]'::jsonb - 0;  
?column?
```

```
-----  
["green", "blue"]
```

#####“#-”#####/#####
json#####contact##fax#/##

```
mydb=> SELECT '{"name": "James", "contact": {"phone": "01234  
567890", "fax": "01987 543210"}}'::jsonb #-  
'{contact,fax}'::text[];  
?column?
```

```
-----  
{"name": "James", "contact": {"phone": "01234 567890"}}  
(1 row)
```

aliases 1 /

```
mydb=> SELECT '{"name": "James", "aliases": ["Jamie", "The
Jamester", "J Man"]}':::jsonb #- '{aliases,1}':::text[];
?column?
```

```
-----
{"name": "James", "aliases": ["Jamie", "J Man"]}
(1 row)
```

json age

```
mydb=> SELECT '{"name":"francs","age":"31"}':::jsonb ||
'{"age":"32"}':::jsonb;
?column?
```

```
-----
{"age": "32", "name": "francs"}
(1 row)
```

jsonb_set

```
jsonb_set(target jsonb, path text[], new_value jsonb[,
create_missing boolean])
```

target jsonb path new_value
create_missing true
create_missing false

```
mydb=> SELECT
jsonb_set('{ "name": "francs", "age": "31" }'::jsonb, '{age}', '"32"
'::jsonb, false);
      jsonb_set
```

```
-----
      {"age": "32", "name": "francs"}
(1 row)
```

```
mydb=> SELECT
jsonb_set('{ "name": "francs", "age": "31" }'::jsonb, '{sex}', '"ma
le"'::jsonb, true);
      jsonb_set
```

```
-----
      {"age": "31", "sex": "male", "name": "francs"}
(1 row)
```

3.9 数据类型

PostgreSQL数据类型转换函数
PostgreSQL数据类型转换函数
CAST数据类型转换函数

3.9.1 数据类型转换函数

PostgreSQL数据类型转换函数3-8

3-8 数据类型转换函数

函 数	返 回 类 型	描 述	示 例
to_char(timestamp, text)	text	把时间戳转换成字符串	to_char(current_timestamp, 'HH12:MI:SS')
to_char(interval, text)	text	把间隔转换成字符串	to_char(interval '15h 2m 12s', 'HH24:MI:SS')
to_char(int, text)	text	把整数转换成字符串	to_char(125, '999')
to_char(numeric, text)	text	把数字转换成字符串	to_char(-125.8, '999D99S')
to_date(text, text)	date	把字符串转换成日期	to_date('05 Dec 2000', 'DD Mon YYYY')
to_number(text, text)	numeric	把字符串转换成数字	to_number('12,454.8-', '99G999D9S')
to_timestamp(text, text)	timestamp with time zone	把字符串转换成时间戳	to_timestamp('05 Dec 2000', 'DD Mon YYYY')

3.9.2 CAST数据类型转换

varchar数据类型text数据类型

```
mydb=> SELECT CAST(vchar'123' as text);
      text
-----
      123
(1 row)
```

varchar → int4

```
mydb=> SELECT CAST(vchar'123' as int4);
      int4
-----
      123
```

3.9.3 数据类型转换

int4 → numeric

```
mydb=> SELECT 1::int4, 3/2::numeric;
 int4 |      ?column?
-----+-----
      1 | 1.5000000000000000
(1 row)
```

SQL 数据库系统元数据表
pg_class 表中的 OID 列

```
mydb=> SELECT oid, relname FROM pg_class WHERE
      relname='test_json1';
      oid | relname
-----+-----
```

```
16509 | test_json1
(1 row)
```

```
test_json1 OID pg_attribute
attrelid OID
```

```
mydb=> SELECT attname FROM pg_attribute WHERE
attrelid='16509' AND attnum >0;
```

```
attname
```

```
-----
```

```
id
```

```
name
```

```
(2 rows)
```

```
mydb=> SELECT attname
FROM pg_attribute
```

```
WHERE attrelid='test_json1'::regclass AND attnum >0;
```

```
attname
```

```
-----
```

```
id
```

```
name
```

```
(2 rows)
```



```
pg_class
```

```
PostgreSQL
```

```
OID pg_class
```

pg_classOIDpg_attribute
pg_attribute.attrelidOID
pgclass.oid

3.10 数据类型

PostgreSQL 数据类型分为基本数据类型和复合数据类型。基本数据类型包括整数、浮点数、字符串、日期、时间、布尔值等。复合数据类型包括数组、JSON、JSONB、XML 等。PostgreSQL 还支持用户自定义数据类型。有关 PostgreSQL 数据类型的详细文档，请访问 <https://www.postgresql.org/docs/10/static/datatype.html>。PostgreSQL 数据类型与 DBA 密切相关。

4 SQL

PostgreSQL SQL
WITH RETURNING
UPSERT

4.1 WITH

WITHはPostgreSQLのSQLの拡張機能で、CTE(Common Table Expressions)と呼ばれる。WITHは、SQLのクエリの中で、一度だけ実行されるクエリを、何度も繰り返し使用できるようにする。これは、SQLのクエリを、より読みやすく、理解しやすくするための機能である。

4.1.1 CTE

CTEはWITHを使用して定義される。

```
WITH t as (  
    SELECT generate_series(1,3)  
)  
SELECT * FROM t;
```

実行結果

```
generate_series  
-----  
1  
2  
3  
(3 rows)
```

CTEは、SQLのクエリの中で、一度だけ実行されるクエリを、何度も繰り返し使用できるようにする。これは、SQLのクエリを、より読みやすく、理解しやすくするための機能である。

CTE를 사용하여 SQL

CTE를 사용하여 SQL을 작성할 때, CTE를 사용하여 데이터를 먼저 처리하고, 그 결과를 기반으로 다른 쿼리를 작성할 수 있습니다. CTE를 사용하여 데이터를 먼저 처리하고, 그 결과를 기반으로 다른 쿼리를 작성할 수 있습니다. CTE를 사용하여 데이터를 먼저 처리하고, 그 결과를 기반으로 다른 쿼리를 작성할 수 있습니다.

```
WITH regional_sales AS (  
    SELECT region, SUM(amount) AS total_sales  
    FROM orders  
    GROUP BY region  
) , top_regions AS (  
    SELECT region  
    FROM regional_sales  
    WHERE total_sales > (SELECT SUM(total_sales)/10 FROM  
regional_sales)  
)  
SELECT region,  
    product,  
    SUM(quantity) AS product_units,  
    SUM(amount) AS product_sales  
FROM orders  
WHERE region IN (SELECT region FROM top_regions)  
GROUP BY region, product;
```

CTE를 사용하여 regional_sales와 top_regions을 정의하고, regional_sales에서 top_regions에 포함된 지역을 필터링하여 orders에서 해당 지역의 판매량을 집계합니다.

4.1.2 CTE를 사용하여

WITH RECURSIVE

RECURSIVE

RECURSIVE

```
WITH recursive t (x) as (  
  SELECT 1  
  UNION  
  SELECT x + 1  
  FROM t  
  WHERE x < 5  
)  
SELECT sum(x) FROM t;
```

sum

```
-----  
      15  
(1 row)
```

SELECT x=1 union 1 SELECT x=5
SELECT x=

PostgreSQL

```
id name fatherid  
1  0  
2  1  
3  1  
4  2
```

```
5 00 2
6 00 3
7 000 4
8 000 4
```

00000id000000000000000id=700000000
000000000id=500000000000000000000000
000000000000000PostgreSQLWITH00000
000000000000000000000000

```
CREATE TABLE test_area(id int4,name varchar(32),fatherid
int4);
```

```
INSERT INTO test_area VALUES (1, '00' ,0);
INSERT INTO test_area VALUES (2, '00' ,1);
INSERT INTO test_area VALUES (3, '00' ,1);
INSERT INTO test_area VALUES (4, '00' ,2);
INSERT INTO test_area VALUES (5, '00' ,2);
INSERT INTO test_area VALUES (6, '00' ,3);
INSERT INTO test_area VALUES (7, '000' ,4);
INSERT INTO test_area VALUES (8, '000' ,4);
```

00PostgreSQLWITH0000ID0700000000
0000000000

```
WITH RECURSIVE r AS (
    SELECT * FROM test_area WHERE id = 7
    UNION ALL
    SELECT test_area.* FROM test_area, r WHERE
test_area.id = r.fatherid
)
SELECT * FROM r ORDER BY id;
```

```
id | name | fatherid
---+-----+-----
  1 |  00  |         0
  2 |  00  |         1
  4 |  00  |         2
  7 |  000 |         4
(4 rows)
```

```
mydb=> WITH RECURSIVE r AS (
    SELECT * FROM test_area WHERE id = 7
    UNION ALL
    SELECT test_area.* FROM test_area, r WHERE
test_area.id = r.fatherid
)
SELECT string_agg(name,'') FROM ( SELECT name FROM r ORDER
BY id) n;
    string_agg
-----
□□□□□□□□
```

```
mydb=> WITH RECURSIVE r AS (
        SELECT * FROM test_area WHERE id = 4
```

```

        UNION ALL
        SELECT test_area.* FROM test_area, r WHERE
test_area.fatherid = r.id
    )
    SELECT * FROM r ORDER BY id;
id | name | fatherid
---+-----+-----
 4 |   |      2
 7 |   |      4
 8 |   |      4
(3 rows)

```

CTE

CTE SQL SQL SQL

CTE

4.2 ☐☐☐☐

PostgreSQL

4.2.1 INSERT INTO...SELECT...

[illegible]

```
INSERT INTO table_name SELECT...FROM source_table
```

```

000000000000user_ini00000000user_ini000
000000000000

```

```
mydb=> CREATE TABLE tbl_batch1(user_id int8,user_name text);
CREATE TABLE
```

```
mydb=> INSERT INTO tbl_batch1(user_id,user_name)
SELECT user_id,user_name FROM user_ini;
INSERT 0 1000000
```

```

        user_ini user_id user_name
tbl_batch1
where

```

```
mydb=> CREATE TABLE tbl_batch2 (id int4,info text);
CREATE TABLE
```

```
mydb=> INSERT INTO tbl_batch2(id,info)
SELECT generate_series(1,5),'batch2';
INSERT 0 5
```

SELECT
PostgreSQL

4.2.2 INSERT INTO VALUES...

PostgreSQL INSERT
VALUES

```
mydb=> CREATE TABLE tbl_batch3(id int4,info text);
CREATE TABLE

mydb=> INSERT INTO tbl_batch3(id,info) VALUES (1,'a'),
(2,'b'),(3,'c');
INSERT 0 3
```

```
mydb=> SELECT * FROM tbl_batch3;
 id | info
-----+-----
  1 | a
  2 | b
```


3 | c
(3 rows)

SQLデータベースはSQL
データベースのWAL (Write-
Ahead Logging) を利用して
PostgreSQLデータベース

4.2.3 COPY \COPY

2.2.3 psqlデータベースCOPY
\copy copy \copy
INSERT PostgreSQL COPY
COPY
4 CPU 8GB

```
mydb=> CREATE TABLE tbl_batch4(  
id int4,  
info text,  
create_time timestamp(6) with time zone default  
clock_timestamp());  
CREATE TABLE
```

```
mydb=> INSERT INTO tbl_batch4(id,info) SELECT n,n||'_batch4'  
FROM generate_series(1,10000000) n;  
INSERT 0 10000000
```

INSERT
[]

```
[postgres@pghost1 ~]$ psql mydb postgres
psql (10.0)
Type "help" for help.
```

```
mydb=# \timing
Timing is on.
```

```
mydb=# COPY pguser.tbl_batch4 TO '/home/pg10/tbl_batch4.txt';
COPY 10000000
Time: 6575.787 ms (00:06.576)
```

6575 tbl_batch4
tbl_batch4.txt

```
mydb=# TRUNCATE TABLE pguser.tbl_batch4;
TRUNCATE TABLE
mydb=# COPY pguser.tbl_batch4 FROM
'/home/pg10/tbl_batch4.txt';
COPY 10000000
Time: 15663.834 ms (00:15.664)
```

COPY 15663

4.3 RETURNING

PostgreSQL RETURNING DML
INSERT RETURNING
UPDATE RETURNING
DELETE RETURNING
SQL

4.3.1 RETURNING

INSERT RETURNING

```
mydb=> CREATE TABLE test_r1(id serial,flag char(1));
CREATE TABLE

mydb=> INSERT INTO test_r1(flag) VALUES ('a') RETURNING *;
   id | flag
-----+-----
    1 | a
(1 row)
INSERT 0 1
```

RETURNING*
RETURNING id

```
mydb=> INSERT INTO test_r1(flag) VALUES ('b') RETURNING id;
      id
-----
        2
(1 row)
INSERT 0 1
```

4.3.2 RETURNING

UPDATE RETURNING UPDATE

```
mydb=> SELECT * FROM test_r1 WHERE id=1;
      id | flag
-----+-----
        1 | a
(1 row)

mydb=> UPDATE test_r1 SET flag='p' WHERE id=1 RETURNING *;
      id | flag
-----+-----
        1 | p
(1 row)
UPDATE 1
```

4.3.3 RETURNING

DELETE RETURNING

```
mydb=> DELETE FROM test_r1 WHERE id=2 RETURNING *;
```

```
  id | flag
```

```
-----+-----
```

```
    2 | b
```

```
(1 row)
```

```
DELETE 1
```

4.4 UPSERT

PostgreSQL UPSERT INSERT...ON
CONFLICT UPDATE
PostgreSQL UPSERT

4.4.1 UPSERT

UPSERT

```
mydb=> CREATE TABLE user_logins(user_name text primary key,  
login_cnt int4,  
last_login_time timestamp(0) without time zone);  
CREATE TABLE
```

```
mydb=> INSERT INTO user_logins(user_name,login_cnt) VALUES  
('francs',1);  
INSERT 0 1
```

user_logins user_name

```
mydb=> INSERT INTO user_logins(user_name,login_cnt)  
VALUES ('matiler',1),('francs',1);
```

ERROR: duplicate key value violates unique constraint
"user_logins_pkey"
DETAIL: Key (user_name)=(francs) already exists.

SQL pour matiler et francs
PostgreSQL UPSERT

```
mydb=> INSERT INTO user_logins(user_name,login_cnt)
VALUES ('matiler',1),('francs',1)
ON CONFLICT(user_name)
DO UPDATE SET
login_cnt=user_logins.login_cnt+EXCLUDED.login_cnt,last_login_time=now();
INSERT 0 2
```

INSERT pour user_logins
login_cnt 1
last_login_time
ON CONFLICT
user_name DO
UPDATE SET
SET user_loins EXCLUDED
user_loins
EXCLUDED user_login

```
mydb=> SELECT * FROM user_logins ;
user_name | login_cnt | last_login_time
```

user_name	login_cnt	last_login_time
matiler	1	
francs	2	2017-08-08 15:23:13

(2 rows)

francs login_cnt
last_login_time matiler

DO NOTHING

```
mydb=> INSERT INTO user_logins(user_name,login_cnt)
VALUES ('tutu',1),('francs',1)
ON CONFLICT(user_name) DO NOTHING;
INSERT 0 1
```

tutu
francs

```
mydb=> SELECT * FROM user_logins ;
 user_name | login_cnt | last_login_time
-----+-----+-----
 matiler   |         1 |
 francs    |         2 | 2017-08-08 15:23:13
 tutu      |         1 |
(3 rows)
```

4.4.2 UPSERT

PostgreSQL UPSERT

```
INSERT INTO table_name [ AS alias ] [ ( column_name [, ...]
) ]
    [ ON CONFLICT [ conflict_target ] conflict_action ]
```

where conflict_target can be one of:

```
( { index_column_name | ( index_expression ) } [ COLLATE
collation ] [ opclass ] [, ...] ) [ WHERE index_predicate ]
ON CONSTRAINT constraint_name
```

and conflict_action is one of:

```
DO NOTHING
DO UPDATE SET { column_name = { expression | DEFAULT } |
                ( column_name [, ...] ) = [ ROW ] ( {
expression | DEFAULT } [, ...] ) |
                ( column_name [, ...] ) = ( sub-SELECT )
                } [, ...]
    [ WHERE condition ]
```

```

[ON
CONFLICT[conflict_target]conflict_action]
conflict_target
conflict_actionDO
NOTHINGUPDATE
```

4.5 随机

TABLESAMPLE 是 PostgreSQL 9.5 引入的。在 PostgreSQL 9.5 之前，使用 ORDER BY random() 来随机选择数据。

```
mydb=> SELECT * FROM user_ini ORDER BY random() LIMIT 1;
      id      | user_id | user_name |      create_time
-----+-----+-----+-----
      500449   |  768810 | 2TY6P4    | 2017-08-05
15:59:32.294761+08
(1 row)
```

```
mydb=> SELECT * FROM user_ini ORDER BY random() LIMIT 1;
      id      | user_id | user_name |      create_time
-----+-----+-----+-----
      324823   |  740720 | 07SKCU    | 2017-08-05
15:59:29.913984+08
(1 row)
```

查询计划

```
mydb=> EXPLAIN ANALYZE SELECT * FROM user_ini ORDER BY
random() LIMIT 1;
                                QUERY PLAN
-----

```

```

Limit (cost=25599.98..25599.98 rows=1 width=35) (actual
time=367.867..367.868 rows=1 loops=1)
  -> Sort (cost=25599.98..28175.12 rows=1030056
width=35) (actual time= 367.866..367.866 rows=1 loops=1)
    Sort Key: (random())
    Sort Method: top-N heapsort  Memory: 25kB
    -> Seq Scan on user_ini (cost=0.00..20449.70
rows=1030056 width=35) (actual time=0.012..159.569
rows=1000000 loops=1)
      Planning time: 0.083 ms
      Execution time: 367.909 ms
(7 rows)

```

user_ini 100 100 SQL
 367ms

9.5 PostgreSQL TABLESAMPLE

```

SELECT ...
FROM table_name
TABLESAMPLE sampling_method ( argument [, ...] ) [
REPEATABLE ( seed ) ]

```

sampling_method
 SYSTEM BERNOLLI
 argument



explain analyze
 SQL SQL Planning time

SQL Execution time SQL

4.5.1 SYSTEM

SYSTEM
SYSTEM
SYSTEM

test_sample150

```
mydb=> CREATE TABLE test_sample(id int4,message text,
create_time timestamp(6) without time zone default
clock_timestamp());
CREATE TABLE
```

```
mydb=> INSERT INTO test_sample(id,message)
SELECT n, md5(random()::text) FROM
generate_series(1,1500000) n;
INSERT 0 1500000
```

```
mydb=> SELECT * FROM test_sample LIMIT 1;
   id | message |
create_time
-----+-----+-----
1 | 58f2506410be948963d6d9adf4b4e0c2 | 2017-08-08
21:17:20.984481
(1 row)
```

0.01%
 $1500000 \times 0.01\% = 150$ SQL

```
EXPLAIN ANALYZE SELECT * FROM test_sample TABLESAMPLE
SYSTEM(0.01);
```

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_sample TABLESAMPLE
SYSTEM(0.01);
```

QUERY PLAN

```
-----
Sample Scan on test_sample (cost=0.00..3.50 rows=150
width=45) (actual time=0.099..0.146 rows=107 loops=1)
  Sampling: system ('0.01'::real)
  Planning time: 0.053 ms
  Execution time: 0.166 ms
(4 rows)
```

Sample Scan
SYSTEM 0.166
150 107 107

```
mydb=> SELECT relname,relpages FROM pg_class WHERE
relname='test_sample';
```

```
      relname | relpages
-----+-----
test_sample |      14019
(1 row)
```

test_sample14019
1000000/14019=107

ctid

```
mydb=> SELECT ctid,* FROM test_sample TABLESAMPLE
SYSTEM(0.01);
   ctid   |   id   |          message          |
create_time
-----+-----+-----+-----
(5640,1)  | 603481 | 385484b3452b245e46388d71ce4ea928 |
2017-08-08 21:17:23.32394
(5640,2)  | 603482 | e09c526118f1d4b3c391d59ae915c4e8 |
2017-08-08 21:17:23.323964
...
(5640,107) | 603587 | c33875a052f4ca63c4b38c649fb6bcc3 |
2017-08-08 21:17:23.324336
(107 rows)
```

ctid
107
5640
0.01

```
mydb=> SELECT count(*) FROM test_sample TABLESAMPLE
SYSTEM(0.01);
   count
-----
      214
(1 row)
```

```
mydb=> SELECT count(*) FROM test_sample TABLESAMPLE
SYSTEM(0.01);
   count
```

107
(1 row)

214107107
0.011500000×0.01%=150
150SYSTEM

4.5.2 BERNOLLI

BERNOULLI
BERNOULLI
BERNOULLI
SYSTEM
SYSTEM
BERNOULLI
test_sample

BERNOULLI0.01

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_sample TABLESAMPLE  
BERNOULLI (0.01);
```

QUERY PLAN

Sample Scan on test_sample (cost=0.00..14020.50
rows=150 width=45) (actual time=0.025..22.541 rows=152
loops=1)

Sampling: bernoulli ('0.01'::real)

Planning time: 0.063 ms
Execution time: 22.569 ms
(4 rows)

Sample Scan
BERNOULLI150152
1501000000×0.01%
22.569SYSTEM0.166
136

```
mydb=> SELECT count(*) FROM test_sample TABLESAMPLE
BERNOULLI(0.01);
count
-----
151
(1 row)
```

```
mydb=> SELECT count(*) FROM test_sample TABLESAMPLE
BERNOULLI(0.01);
count
-----
147
(1 row)
```

BERNOULLI
SYSTEM
SYSTEM

TABLESAMPLE BERNOLLI
ctid

```
mydb=> SELECT ctid,id,message
        FROM test_sample TABLESAMPLE BERNOLLI(0.01) LIMIT
3;
```

ctid	id	message
(55,30)	5915	f3803f234f6cf6cdd276d9d027487582
(240,23)	25703	c04af69ac76f6465832e0cd87939a1af
(318,3)	34029	dd35438b24980d1a8ed2d3f5edd5ca1c

ctid55
240318BERNOLLI
SYSTEM

SYSTEMBERNOLLI
SYSTEM
GB
BERNOLLI
SYSTEM
SYSTEM

4.6 聚合函数

聚合函数用于对一组值进行计算并返回一个单一的值。常见的聚合函数包括 avg、sum、min、max、count 等。PostgreSQL 提供了丰富的聚合函数，可以满足各种需求。

聚合函数通常与 GROUP BY 子句一起使用，以对数据进行分组并计算每个组的聚合值。聚合函数也可以与 HAVING 子句一起使用，以对聚合结果进行过滤。

聚合函数	返回类型
avg	numeric
sum	numeric
min	numeric
max	numeric
count	integer

聚合函数可以用于计算数据的平均值、总和、最小值、最大值和计数等。

聚合函数	聚合函数
avg	avg
sum	sum

聚合函数可以用于计算 SQL 语句的结果。

4.6.1 string_agg

string_agg 聚合函数用于将一组字符串聚合成一个单一的字符串。

```
string_agg(expression, delimiter)
```

string_agg aggregates values of expression into a single string or bytea value, separated by the delimiter. The expression must be of type text or bytea. The delimiter must be of type text or bytea. The result is of type text or bytea, depending on the type of the expression.

Example:

```
CREATE TABLE city (country character varying(64),city
character varying(64));
INSERT INTO city VALUES ('US','New York');
INSERT INTO city VALUES ('US','Los Angeles');
INSERT INTO city VALUES ('US','Chicago');
INSERT INTO city VALUES ('US','Houston');
INSERT INTO city VALUES ('US','Phoenix');
```

Query:

```
mydb=> SELECT * FROM city;
      country | city
```

```
-----+-----
      US      | New York
      US      | Los Angeles
      US      | Chicago
      US      | Houston
      US      | Phoenix
```

```
(5 rows)
```

city

```
mydb=> SELECT string_agg(city,',') FROM city;
      string_agg
```

```
-----
    00,00,00,00,00
(1 row)
```

string_agg

SQL

```
mydb=> SELECT country,string_agg(city,',') FROM city GROUP
BY country;
```

```
   country | string_agg
-----+-----
    00    | 00,00
    00    | 00,00,00
(2 rows)
```

4.6.2 array_agg

array_agg

string_agg

array_agg

```
array_agg(expression)  --
```

array_agg

```
mydb=> SELECT country,array_agg(city) FROM city GROUP BY
country;
```

country	array_agg
US	{CA,PA}
US	{CA,PA,WA}

array_agg

array_agg

array_agg(expression) --

array_agg

```
mydb=> CREATE TABLE test_array3(id int4[]);
CREATE TABLE
mydb=> INSERT INTO test_array3(id) VALUES (array[1,2,3]);
INSERT 0 1
mydb=> INSERT INTO test_array3(id) VALUES (array[4,5,6]);
INSERT 0 1
```

array_agg

```
mydb=> SELECT * FROM test_array3;
      id
-----
    {1,2,3}
    {4,5,6}
(2 rows)
```

array_agg

```
mydb=> SELECT  array_agg(id) FROM test_array3;
      array_agg
-----
 {{1,2,3},{4,5,6}}
(1 row)
```

array_agg
array_to_string

```
mydb=> SELECT array_to_string( array_agg(id),',') FROM
test_array3;
      array_to_string
-----
    1,2,3,4,5,6
(1 row)
```

4.7 窗口函数

窗口函数是SQL中非常重要的一类函数，它可以在不改变表结构的情况下，对表中的数据进行分组、排序、聚合等操作。窗口函数的使用可以大大简化SQL语句，提高查询效率。本章将详细介绍窗口函数的使用方法和注意事项。

4.7.1 窗口函数概述

PostgreSQL窗口函数包括row_num、rank、lag等函数。这些函数可以在查询结果中返回每一行的排名、位置等信息。窗口函数的使用通常与OVER子句结合使用。

窗口函数的基本语法如下：

```
function_name ([expression [, expression ... ]]) [ FILTER (
WHERE filter_clause ) ] OVER ( window_definition )
```

其中window_definition的定义如下：

```
[ existing_window_name ]
[ PARTITION BY expression [, ...] ]
[ ORDER BY expression [ ASC | DESC | USING operator ] [
NULLS { FIRST | LAST } ] [, ...] ]
[ frame_clause ]
```

□□□□□

·OVER□□□□□□□□□□

·PARTITION BY□□□□□□□□□□□□□□□□□□
□□□□□□□□□□

·ORDER BY□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□

4.7.2 avg□□OVER□□

□□□□□□OVER□□□□□□□□□□□□□□□□□□□□
□□□□□□□□avg□□□□□□□□OVER□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□

```
CREATE TABLE score ( id serial primary key,
                        subject character varying(32),
                        stu_name character varying(32),
                        score numeric(3,0) );
```

```
INSERT INTO score ( subject,stu_name,score ) VALUES
('Chinese','francs',70);
```

```
INSERT INTO score ( subject,stu_name,score ) VALUES
('Chinese','matiler',70);
```

```
INSERT INTO score ( subject,stu_name,score) VALUES
('Chinese','tutu',80);
```

```
INSERT INTO score ( subject,stu_name,score ) VALUES
```


subject	stu_name	score	avg
Chinese	francs	70	73.33333333333333
Chinese	matiler	70	73.33333333333333
Chinese	tutu	80	73.33333333333333
English	matiler	75	75.00000000000000
English	francs	90	75.00000000000000
English	tutu	60	75.00000000000000
Math	francs	80	81.33333333333333
Math	matiler	99	81.33333333333333
Math	tutu	65	81.33333333333333

(9 rows)

PARTITION BY subject

4.7.3 row_number

row_number

```
mydb=> SELECT row_number() OVER (partition by subject ORDER BY score desc),* FROM score;
```

row_number	id	subject	stu_name	score
1	3	Chinese	tutu	80
2	1	Chinese	francs	70
3	2	Chinese	matiler	70
1	5	English	francs	90
2	4	English	matiler	75
3	6	English	tutu	60
1	8	Math	matiler	99
2	7	Math	francs	80
3	9	Math	tutu	65

(9 rows)

row_number() over (order by id) as rownum ,*
 FROM score;
 rownum | id | subject | stu_name | score
 -----+-----+-----+-----+-----
 1 | 1 | Chinese | francs | 70
 2 | 2 | Chinese | matiler | 70
 3 | 3 | Chinese | tutu | 80
 4 | 4 | English | matiler | 75
 5 | 5 | English | francs | 90
 6 | 6 | English | tutu | 60
 7 | 7 | Math | francs | 80
 8 | 8 | Math | matiler | 99
 9 | 9 | Math | tutu | 65

(9 rows)

4.7.4 rank()

rank() over (partition by subject order by score),* FROM score;
 rank | id | subject | stu_name | score
 -----+-----+-----+-----+-----
 1 | 2 | Chinese | matiler | 70
 1 | 1 | Chinese | francs | 70
 3 | 3 | Chinese | tutu | 80
 1 | 6 | English | tutu | 60
 2 | 4 | English | matiler | 75

mydb=> SELECT rank() OVER(PARTITION BY subject ORDER BY
 score),* FROM score;

rank | id | subject | stu_name | score
 -----+-----+-----+-----+-----
 1 | 2 | Chinese | matiler | 70
 1 | 1 | Chinese | francs | 70
 3 | 3 | Chinese | tutu | 80
 1 | 6 | English | tutu | 60
 2 | 4 | English | matiler | 75

3	5	English	francs	90
1	9	Math	tutu	65
2	7	Math	francs	80
3	8	Math	matiler	99

(9 rows)

查询Chinese科目中分数为70的学生排名为1，分数为80的学生排名为3

4.7.5 dense_rank

dense_rank函数返回的排名与rank函数返回的排名类似，但dense_rank函数返回的排名是连续的，而rank函数返回的排名是不连续的。

```
mydb=> SELECT dense_rank() OVER(PARTITION BY subject ORDER BY score),* FROM score;
```

dense_rank	id	subject	stu_name	score
1	2	Chinese	matiler	70
1	1	Chinese	francs	70
2	3	Chinese	tutu	80
1	6	English	tutu	60
2	4	English	matiler	75
3	5	English	francs	90
1	9	Math	tutu	65
2	7	Math	francs	80
3	8	Math	matiler	99

(9 rows)

查询Chinese科目中排名为1的学生排名为1，排名为2的学生排名为2

4.7.6 lag

lag 関数は、offset 引数で指定した行数前のデータを取得する。
引数は、value、offset、default の 3 つである。

```
lag(value anyelement [, offset integer [, default anyelement]])
```

引数

- ・value: 取得するデータの列名

- ・offset: 取得するデータの行数。省略時は 1。

- ・default: offset が 0 の場合の default 値。省略時は null。

例: score テーブルの id 列のデータを取得する。

```
mydb=> SELECT lag(id,1)OVER(* FROM score;
```

lag	id	subject	stu_name	score
	1	Chinese	francs	70
1	2	Chinese	matiler	70
2	3	Chinese	tutu	80
3	4	English	matiler	75
4	5	English	francs	90
5	6	English	tutu	60
6	7	Math	francs	80

7	8	Math	matiler	99
8	9	Math	tutu	65

(9 rows)

score id

```
mydb=> SELECT lag(id,2,1000)OVER(* FROM score;
```

lag	id	subject	stu_name	score
1000	1	Chinese	francs	70
1000	2	Chinese	matiler	70
1	3	Chinese	tutu	80
2	4	English	matiler	75
3	5	English	francs	90
4	6	English	tutu	60
5	7	Math	francs	80
6	8	Math	matiler	99
7	9	Math	tutu	65

(9 rows)

lag offset

4.7.7 first_value

first_value

score

```
mydb=> SELECT first_value(score) OVER( PARTITION BY subject
),* FROM score;
```

first_value	id	subject	stu_name	score
70	1	Chinese	francs	70
70	2	Chinese	matiler	70
70	3	Chinese	tutu	80
75	4	English	matiler	75
75	5	English	francs	90
75	6	English	tutu	60
80	7	Math	francs	80
80	8	Math	matiler	99
80	9	Math	tutu	65

(9 rows)

first_value
score

```
mydb=> SELECT first_value(score) OVER( PARTITION BY subject
ORDER BY score desc),* FROM score;
```

first_value	id	subject	stu_name	score
80	3	Chinese	tutu	80
80	1	Chinese	francs	70
80	2	Chinese	matiler	70
90	5	English	francs	90
90	4	English	matiler	75
90	6	English	tutu	60
99	8	Math	matiler	99
99	7	Math	francs	80
99	9	Math	tutu	65

(9 rows)

4.7.8 last_value

last_value() over() over() over()

score() over() over() over()

```
mydb=> SELECT last_value(score) OVER( PARTITION BY subject
),* FROM score;
```

last_value	id	subject	stu_name	score
80	1	Chinese	francs	70
80	2	Chinese	matiler	70
80	3	Chinese	tutu	80
60	4	English	matiler	75
60	5	English	francs	90
60	6	English	tutu	60
65	7	Math	francs	80
65	8	Math	matiler	99
65	9	Math	tutu	65

(9 rows)

4.7.9 nth_value()

nth_value() over() over() over()

nth_value(value any, nth integer)

over()

value() over()

nth

score

```
mydb=> SELECT nth_value(score,2) OVER( PARTITION BY subject
),* FROM score;
```

nth_value	id	subject	stu_name	score
70	1	Chinese	francs	70
70	2	Chinese	matiler	70
70	3	Chinese	tutu	80
90	4	English	matiler	75
90	5	English	francs	90
90	6	English	tutu	60
99	7	Math	francs	80
99	8	Math	matiler	99
99	9	Math	tutu	65

(9 rows)

4.7.10

SQL

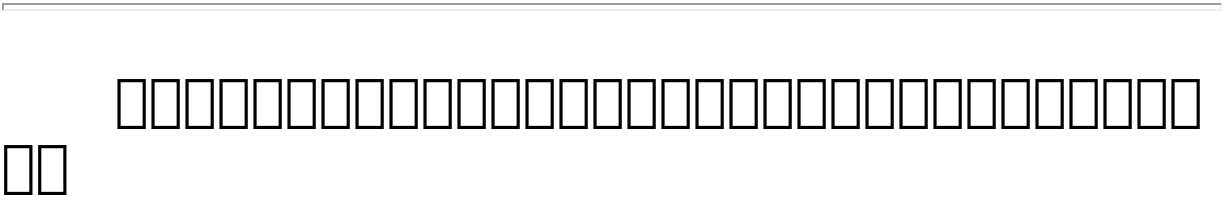
```
SELECT .. FROM .. WINDOW window_name AS ( window_definition
) [, ...]
```

WINDOW window_name
OVER

```
mydb=> SELECT avg(score) OVER(r),sum(score) OVER(r),* FROM
SCORE WINDOW r as (PARTITION BY subject);
```

	avg	sum	id	subject	stu_name	
score						
	-----+-----+-----+-----+-----+-----					

70	73.33333333333333	220	1	Chinese	francs	
70	73.33333333333333	220	2	Chinese	matiler	
80	73.33333333333333	220	3	Chinese	tutu	
75	75.00000000000000	225	4	English	matiler	
90	75.00000000000000	225	5	English	francs	
60	75.00000000000000	225	6	English	tutu	
80	81.33333333333333	244	7	Math	francs	
99	81.33333333333333	244	8	Math	matiler	
65	81.33333333333333	244	9	Math	tutu	
(9 rows)						



4.8 数据库

数据库 PostgreSQL 数据库 SQL 数据库
数据库 SQL 数据库 数据库 数据库 数据库 数据库 数据库 数据库
数据库 数据库 数据库 数据库 数据库 SQL 数据库 数据库
PostgreSQL 数据库 SQL 数据库

データベース

- ・第5回 データベース
- ・第6回 データベース
- ・第7回 データベース
- ・第8回 データベース
- ・第9回 PostgreSQLとNoSQL

5 数据库

PostgreSQL数据库系统是一个开源的、关系型的数据库系统。它支持SQL语言，并且具有强大的事务处理能力。PostgreSQL数据库系统是由PostgreSQL基金会维护的。PostgreSQL数据库系统是一个开源的、关系型的数据库系统。它支持SQL语言，并且具有强大的事务处理能力。PostgreSQL数据库系统是由PostgreSQL基金会维护的。PostgreSQL数据库系统是一个开源的、关系型的数据库系统。它支持SQL语言，并且具有强大的事务处理能力。PostgreSQL数据库系统是由PostgreSQL基金会维护的。

5.1 〇〇〇〇〇〇〇〇〇〇

```

PostgreSQL Database
Cluster
PostgreSQL
"
PostgreSQL

```

5.1.1 五五五五五

```

Database[] User[] Database[] Database
[] User[] 5-1 []
[]

```

```

initdb

```

1. 環境構築

PostgreSQLのインストールと
PGDATAの指定方法について
確認する

```
[postgres@pghost1 ~]$ tree -L 1 -d /pgdata/10/data
/pgdata/10/data
├── base
├── pg_tblspc
├── ...
├── ...
├── ...
├── pg_wal
└── global
```

5-1 インストールと設定

5-1 インストールと設定

目 录	用 途
base	包含每个数据库对应的子目录的子目录
global	包含集群范围的表的子目录, 比如 pg_database
pg_commit_ts	包含事务提交时间戳数据的子目录
pg_xact	包含事务提交状态数据的子目录
pg_dynshmem	包含被动态共享内存子系统所使用文件的子目录
pg_logical	包含用于逻辑复制的状态数据的子目录
pg_multixact	包含多事务状态数据的子目录 (用于共享的行锁)
pg_notify	包含 LISTEN/NOTIFY 状态数据的子目录
pg_repslot	包含复制槽数据的子目录
pg_serial	包含已提交的可序列化事务信息的子目录
pg_snapshots	包含导出的快照的子目录
pg_stat	包含用于统计子系统的永久文件的子目录
pg_stat_tmp	包含用于统计信息子系统临时文件的子目录
pg_subtrans	包含子事务状态数据的子目录
pg_tblspc	包含指向表空间的符号链接的子目录
pg_twophase	用于预备事务状态文件的子目录
pg_wal	保存预写日志
pg_xact	记录事务提交状态数据
文 件	用 途
PG_VERSION	PostgreSQL 主版本号文件
pg_hba.conf	客户端认证控制文件
postgresql.conf	参数文件
postgresql.auto.conf	参数文件, 只保存 ALTER SYSTEM 命令修改的参数
postmaster.opts	记录服务器最后一次启动时使用的命令行参数

2. 数据目录

数据目录 base 是数据库的根目录，它包含所有数据库的数据文件。每个数据库的数据文件都存储在 base 目录下的一个子目录中。每个子目录的名称是该数据库的 OID。

1. OID

PostgreSQLデータベースのOID
OIDはデータベースごとに4バイトのOID
で識別される。OIDはpg_database
のOID

```
SELECT oid,datname FROM pg_database WHERE datname = 'mydb';
oid    | datname
-----+-----
16384  | mydb
(1 row)
```

OIDはpg_classの
OID

```
mydb=# SELECT oid,relname,relkind FROM pg_class WHERE
relname ~ 'tbl';
oid    | relname                | relkind
-----+-----+-----
16385  | tbl_id_seq             | S
16387  | tbl                    | r
16396  | tbl_pkey               | i
3455   | pg_class_tblspc_relfilenode_index | i
(4 rows)
```

2

PostgreSQLデータベースのOID
OIDはデータベースごとに4バイトのOID
で識別される。OIDはpg_database
のOID

pg_default pg_global

```
mydb=# \db
      List of tablespaces
      Name      |  Owner   | Location
-----+-----+-----
 pg_default    | postgres |
 pg_global     | postgres |
(2 rows)
```

pg_global

pg_default base
template0 template1
template1
template1

WAL

postgresql
CREATE TABLESPACE

```
[postgres@pghost1 ~]$ mkdir -p /pgdata/10/mytblspc
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/psql -p 1921 mydb
psql (10.2)
Type "help" for help.
mydb=# CREATE TABLESPACE myspc LOCATION
'/pgdata/10/mytblspc';
CREATE TABLESPACE
mydb=# \db
```

```
          List of tablespaces
   Name      |  Owner   |      Location
-----+-----+-----
myspc        | postgres | /pgdata/10/mytblspc
pg_default   | postgres |
pg_global    | postgres |
(3 rows)
```

```
mydb=# CREATE TABLE t(id SERIAL PRIMARY KEY, ival int)
TABLESPACE myspc;
CREATE TABLE
```

OID
mydb
pg_databaseOID

```

mydb=# SELECT oid,datname FROM pg_database WHERE datname =
'mydb';
      oid | datname
-----+-----
    16384 | mydb
(1 row)

```

```

cd /var/lib/pgsql/data/mydb/OID/16384/
mydb$PGDATA/base/16384/

```

```

[postgres@pghost1 ~]$ ll /pgdata/10/data/base/16384/
-rw----- 1 postgres postgres 16384 Nov 28 21:22 3712
...
...
-rw----- 1 postgres postgres 8192 Nov 28 21:22 3764_vm

```

3

```

PostgreSQL
OID 1GB PostgreSQL
OID.<>
OID.<>
pg_class relfilenode
pg_class OID
relfilenode VACUUM TRUNCATE
relfilenode

```

□□□□□

```
mydb=# SELECT oid,relfilenode FROM pg_class WHERE relname =
'tbl';
      oid | relfilenode
-----+-----
    16387 |          16387
(1 row)
mydb=# \! ls -l /pgdata/10/data/base/16384/16387*
-rw----- 1 postgres postgres 8192 Mar 26 22:22
/pgdata/10/data/base/16384/16387
```

□□□□□tbl□□OID□16387□relfilenode□□
16387□□□□□□□
□"/pgdata/10/data/base/16384/16387"□□□
TRUNCATE□□tbl□□□□□□□□□□□□

```
mydb=# TRUNCATE tbl;
TRUNCATE TABLE
mydb=# CHECKPOINT;
CHECKPOINT
mydb=# \! ls -l /pgdata/10/data/base/16384/16387*
ls: cannot access /pgdata/10/data/base/16384/16387*: No such
file or directory
```

□□□□□□□tbl□□□□□□□
□"/pgdata/10/data/base/16384/16387"□□□
□□□□□tbl□□□□□□□□□□□□

```
postgres@160.40:1922/mydb=# select oid,relfilenode from
pg_class where relname = 'tbl';
      oid | relfilenode
```

```
-----+-----
      16387 |          24591
(1 row)
postgres@160.40:1922/mydb=# \! ls -l
/pgdata/10/data/base/16384/24591*
-rw----- 1 postgres postgres 0 Apr  2 21:24
/pgdata/10/data/base/16384/24591
```

pg_classtbl
"/pgdata/10/data/base/16384/24591"
<relfilenode>.<>

tbl

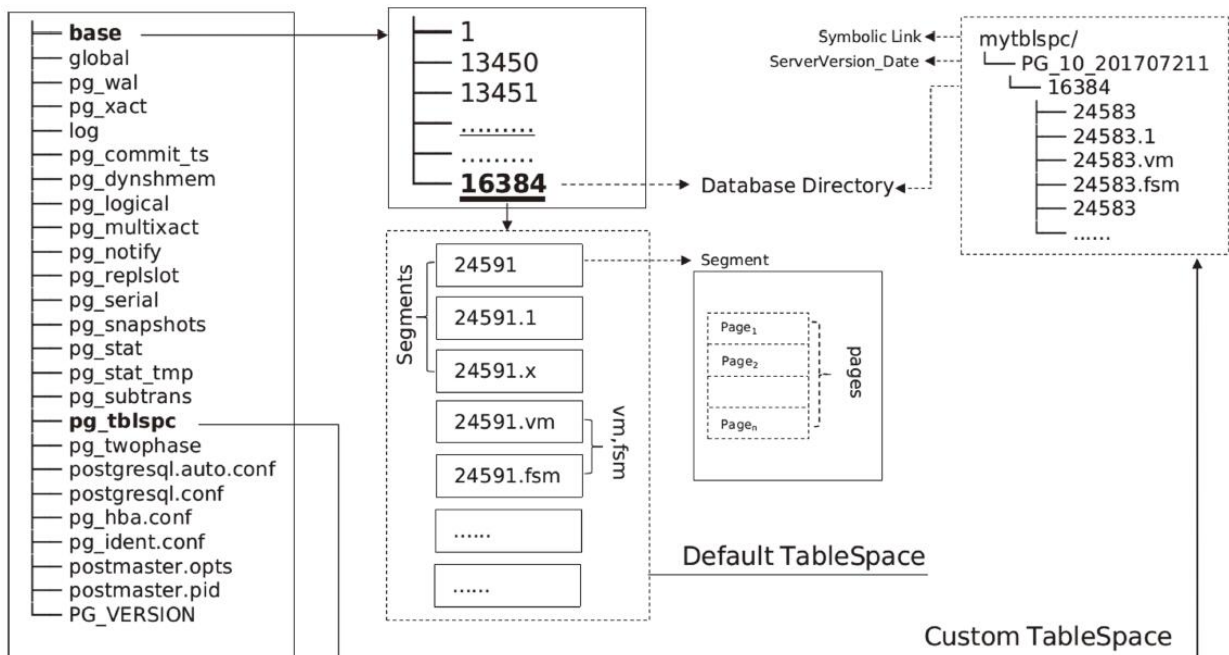
```
mydb=# insert into tbl (ival,description,created_time)
select (random()*(2*10^9)):: integer as
ival,substr('abcdefghijklmnopqrstuvwxyz',1,
(random()*26)::integer) as
description,date(generate_series(now(), now() + '1 week', '1
day')) as created_time from generate_series(1,2000000);
INSERT 0 16000000
```

```
mydb=# SELECT
pg_size_pretty(pg_relation_size('tbl'::regclass));
pg_size_pretty
-----
      1068 MB
(1 row)
```

tbl1068MB
UPDATE

```
/mydb=# \! ls -lh /pgdata/10/data/base/16384/24591*  
-rw----- 1 postgres postgres 1.0G Apr  7 08:44  
/pgdata/10/data/base/16384/24591  
-rw----- 1 postgres postgres 383M Apr  7 08:44  
/pgdata/10/data/base/16384/24591.1  
-rw----- 1 postgres postgres 376K Apr  7 08:44  
/pgdata/10/data/base/16384/24591_fsm  
-rw----- 1 postgres postgres 8.0K Apr  7 08:44  
/pgdata/10/data/base/16384/24591_vm
```

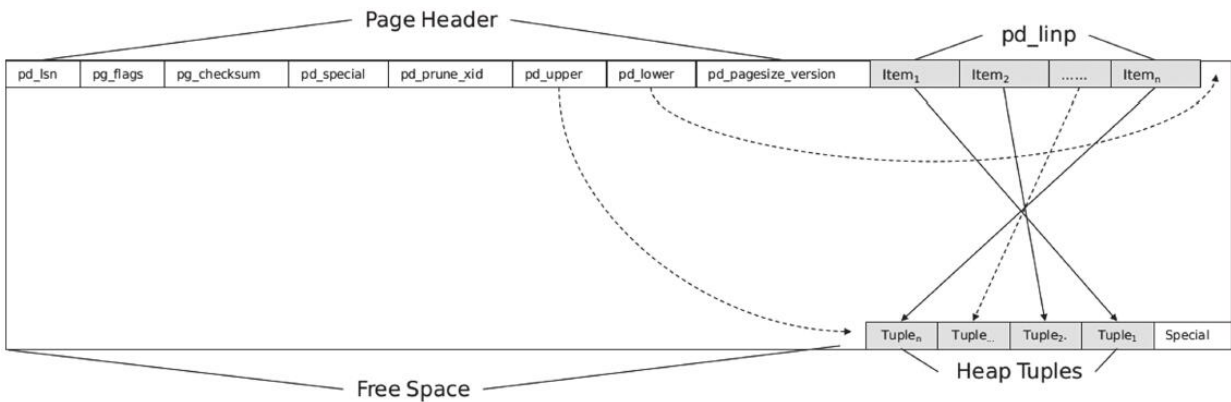
<relfilenode>.<
>tbl1GBtblrelfilenode24591
1GBGB
24591.1_fsm_vm
5-2
PostgreSQL



5-2 PostgreSQLのファイル構造

4. データの保存方法

PostgreSQLはデータをPage単位で保存し、PageはBufferプールに格納される。RelationはTupleの集合であり、Pageは8kBのサイズを持つ。PostgreSQLはBLCKSZパラメータでPageのサイズを指定し、PageはTupleの集合である。I/O操作はPage単位で行われ、BLCKSZパラメータで指定される。



5-3 Page

PageHeader
 PageHeader

·pd_lsn ARIES Recovery Algorithm
 lsn PageLSN xlog
 LSN WAL
 pd_lsn xlogid
 64 LSN 32

·pg_flags

·pd_special

·pd_lower

·pd_upper

·pd_page_size_version PostgreSQL
0000000000000000

·pd_linp[1] 5-3 Item1
Item2...Itemn Tuple

pd_lower pd_upper
pd_lower pd_upper
pd_lower pd_upper
pd_lsn pg_checksum pg_flag

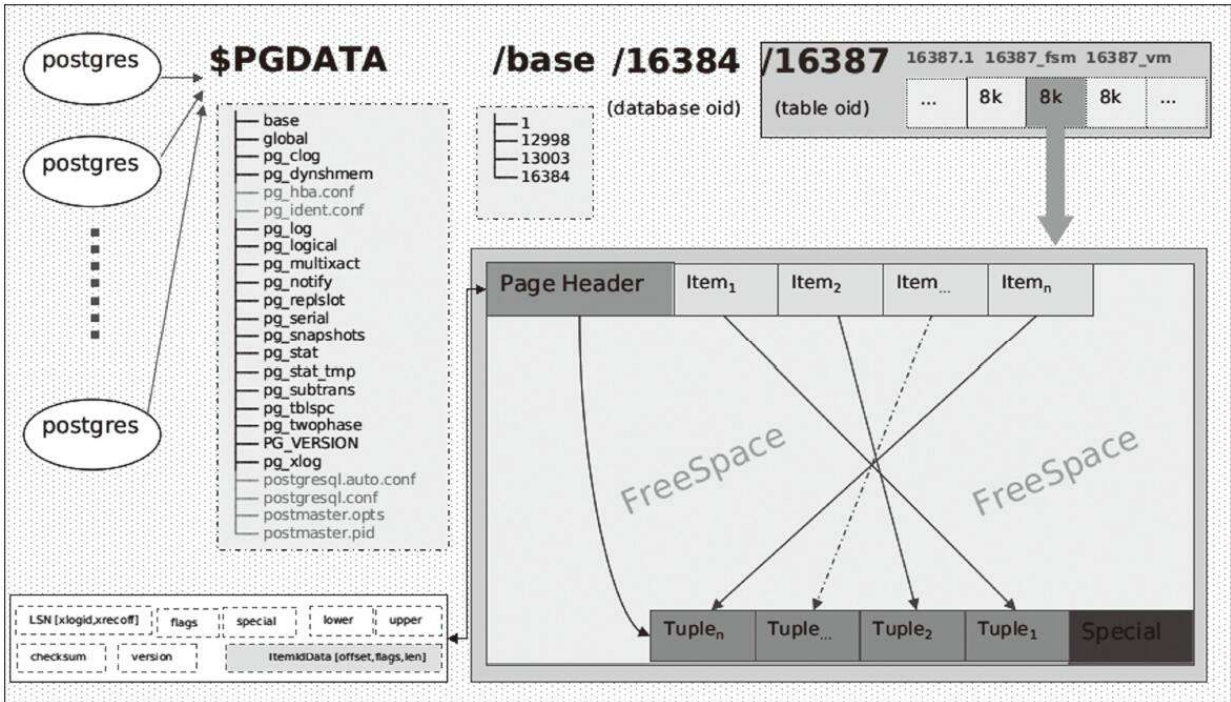
B
B
TID PostgreSQL
TID

Tuple
HeapTupleHeader Tuple 5-4
Tuple OID xmin cmin
HeapTuple Tuple

OID	- object ID of tuple	Header
xmin	- creation transaction ID	
xmax	- destruction transaction ID	
cmin	- creation command ID	
cmax	- destruction command ID	
ctid	- tuple ID (page,item)	
natts	- number of attributes	
infomask	- tuple flags	
hoff	- length of tuple header	
bits	- bit map representing NULLs	
.....		
.....		Values
.....		
.....		
.....		

図5-4 Tupleの構造

図5-5 属性値の格納方法



□5-5 □□□□□□

5.2 进程

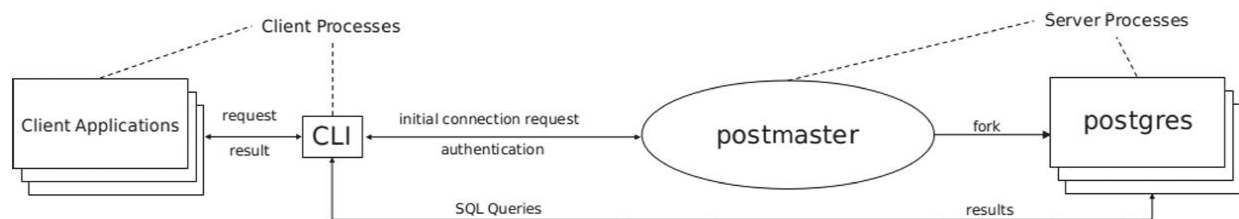
PostgreSQL 进程包括以下/进程：
postmaster 主进程
postgres 子进程
syslogger 系统日志进程
checkpointer 检查点进程
bgwriter 后台写进程
walwriter 写前日志进程

5.2.1 主进程

postmaster 主进程是 PostgreSQL 的入口点，它负责启动其他子进程。

- 启动子进程
- 接收客户端连接
- 管理数据库的 fork 子进程 postgres
- 管理数据库的共享内存
- 管理数据库的锁
- 管理数据库的日志

postgresql.conf
 postmaster fork postgres
 postgres postgres
 postmaster postgres postgres
 postgres 5-6



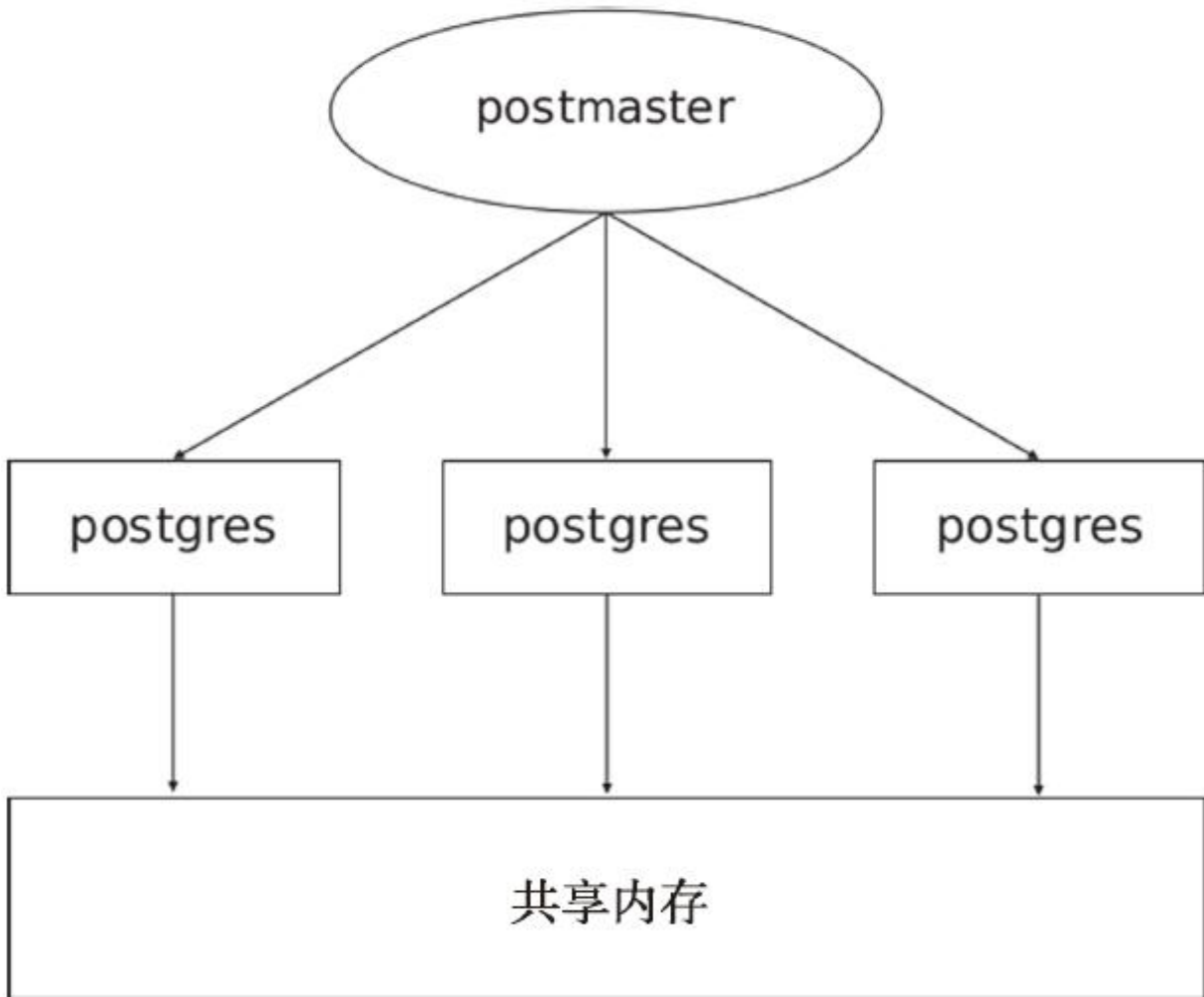
5-6 PostgreSQL architecture

PostgreSQL
 PostgreSQL
 PostgreSQL TCP/IP Unix
 PostgreSQL
 PostgreSQL 5-7

5.2.2

postmaster postgres
 PostgreSQL

·background writer bgwriter
 bgwriter
 postgres



□5-7 □□□□□□□□□□

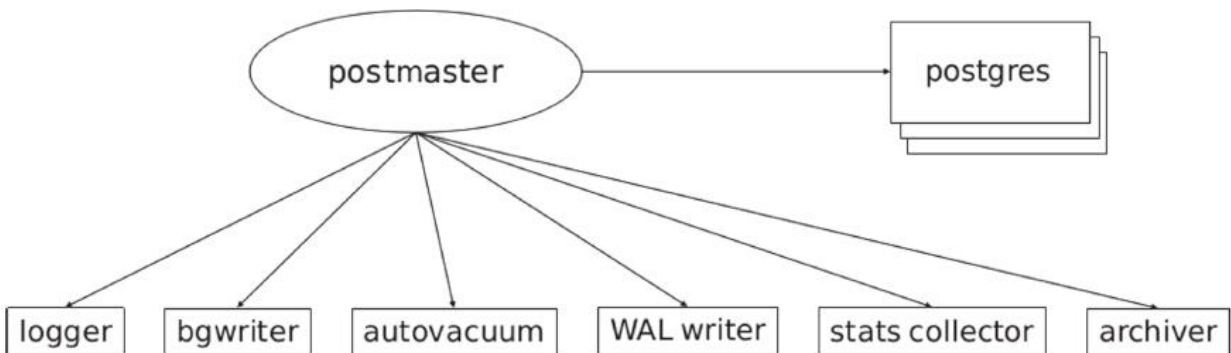
- autovacuum launcher□□□□□□□□□□
- WAL writer□□□□WAL□□□□□WAL□□□□□□□
- statistics collector□□□□□□□□□□
- logging collector□□□□□□□□□□□□□□□□□□

□□

·archiver WAL

·checkpointer

5-8 postmaster



5-8

5.3 準備

PostgreSQLのインストールと起動については、[こちら](#)のページを参照してください。

5.3.1 環境

以下のコマンドを実行して、`fork`、`work_mem`、`maintenance_work_mem`、`temp_buffers`の値を確認してください。

```
・work_mem → ORDER BY、DISTINCTの  
クエリで使用するメモリ
```

```
・maintenance_work_mem →  
VACUUM、REINDEX、CREATE INDEXの  
クエリ
```

```
・temp_buffers → テンポラリテーブルのメモリ
```

5.3.2 設定

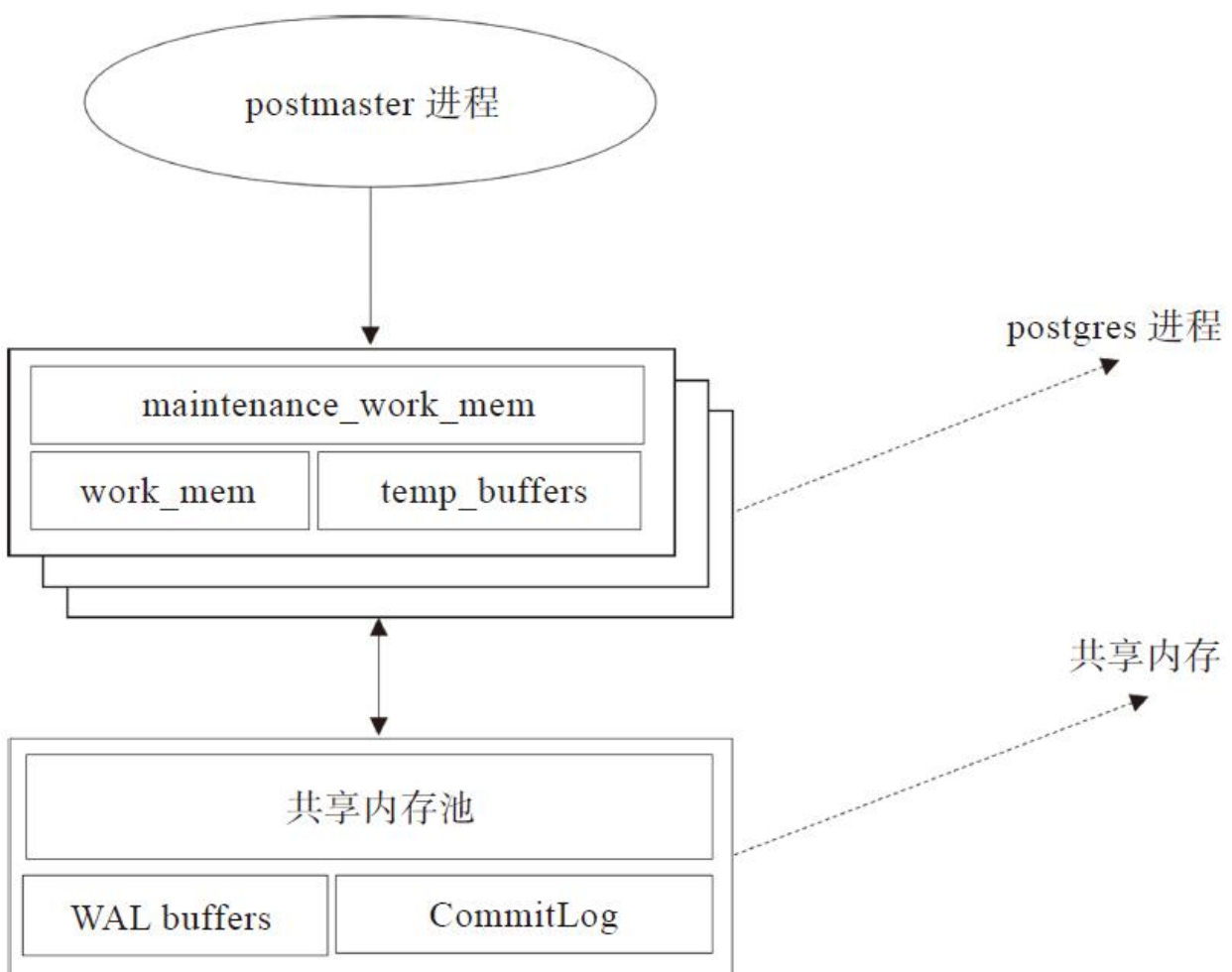
以下のコマンドを実行して、PostgreSQLのインストールディレクトリを確認してください。

·shared buffer pool PostgreSQL 数据库系统
数据库系统共享内存池

·WAL buffer WAL 数据库系统

·CommitLog buffer PostgreSQL
Commit Log 数据库系统
数据库系统

5-9 数据库系统



□5-9 □□□□

5.4 数据库

数据库是存储和管理数据的系统。PostgreSQL 是一个开源的数据库系统，它支持多种数据类型和复杂的查询。数据库的设计和管理对于应用程序的性能和安全性至关重要。PostgreSQL 提供了强大的功能，包括事务处理、并发控制和安全性。它支持多种数据类型，包括整数、浮点数、字符串和日期。PostgreSQL 还支持复杂的查询，包括子查询、聚合函数和连接。数据库的设计和管理需要考虑许多因素，包括性能、安全性和可扩展性。PostgreSQL 是一个成熟且稳定的数据库系统，适用于各种规模的应用程序。它提供了丰富的文档和社区支持，使其成为开发者的首选数据库之一。

第6章 数据库

Oracle数据库Oracle数据库
SELECT UPDATE DELETE
CPU PostgreSQL
9.6 SQL CPU
9.6 9.6
10
index-only
bitmap heap PostgreSQL 10

6.1 数据库配置

PostgreSQL数据库配置参数

1.max_worker_processesinteger

数据库最大并行进程数，默认为8，即最多有8个并行进程同时运行。

2.max_parallel_workersinteger

数据库最大并行进程数，默认为8，即最多有8个并行进程同时运行。
max_worker_processes数据库最大并行进程数
max_parallel_workers_per_gather数据库最大并行进程数

数据库最大并行进程数，默认为8，即最多有8个并行进程同时运行。
max_parallel_workers_per_gather数据库最大并行进程数

3.max_parallel_workers_per_gatherinteger

数据库最大并行进程数，默认为2，即最多有2个并行进程同时运行。
max_worker_processes数据库最大并行进程数
max_parallel_workers_per_gather数据库最大并行进程数

データベースはOLTPデータベースとして
構築する

postgresql.confを編集する

```
max_worker_processes = 16
max_parallel_workers_per_gather = 4      # taken from
max_parallel_workers
max_parallel_workers = 8
parallel_tuple_cost = 0.1
parallel_setup_cost = 1000.0
min_parallel_table_scan_size = 8MB
min_parallel_index_scan_size = 512kB
force_parallel_mode = off
```

4CPU 8GB
環境



パラメータ設定

```
max_parallel_workers_per_gatherを16に
データベースのSQLを実行する際に
max_parallel_workersを
max_worker_processesを16にSQLを実行
する際に
max_worker_processesを2に
max_parallel_workers_per_gatherを4に
データベースのSQLを実行する際に2に
min_parallel_table_scan_sizeを8MB
```

Workers
Planned Worker
Launched

6.2 索引

PostgreSQL 索引类型包括 B-tree、hash、index-only、bitmap heap 等。index-only 索引只包含索引键，不包含数据行指针。bitmap heap 索引用于堆表，通过 bitmap 来定位数据行。

6.2.1 索引类型

索引类型包括 sequential scan、index scan、bitmap scan 等。sequential scan 是顺序扫描，index scan 是索引扫描，bitmap scan 是 bitmap 扫描。CPU 和 IO 是索引性能的关键因素。OLTP 系统通常使用索引，而 OLAP 系统通常使用 sequential scan。

创建索引的 SQL 语句如下：

```
CREATE TABLE test_big1(  
  id int4,  
  name character varying(32),  
  create_time timestamp without time zone default  
  clock_timestamp());  
  
INSERT INTO test_big1(id,name)  
SELECT n, n|| '_test' FROM generate_series(1,50000000) n ;
```

索引的创建和删除操作如下：

```
mydb=> EXPLAIN SELECT * FROM test_big1 WHERE name='1_test';
      QUERY PLAN
```

```
-----
--
      Seq Scan on test_big1  (cost=0.00..991664.00 rows=1
width=25)
      Filter: ((name)::text = '1_test'::text)
(2 rows)
```

Seq Scan on test_big1
test_big1
PostgreSQL 9.6
CPU

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_big1 WHERE
name='1_test';
```

```
      QUERY PLAN
```

```
-----
-----
      Gather  (cost=1000.00..523914.10 rows=1 width=25) (actual
time=0.440..1362.675 rows=1 loops=1)
        Workers Planned: 4
        Workers Launched: 4
        -> Parallel Seq Scan on test_big1  (cost=0.00..522914.00
rows=1 width=25) (actual
time=1083.280..1355.685 rows=0 loops=5)
          Filter: ((name)::text = '1_test'::text)
          Rows Removed by Filter: 10000000
        Planning time: 0.085 ms
        Execution time: 1367.248 ms
(8 rows)
```

Workers Planned
Worker Launched

```
Workers PlannedWorker Launched
4Parallel Seq Scan on test_big1
Planning timeExecution
timeSQL4SQL
1367SQL
max_parallel_workers_per_gather4
00000
```

```
mydb=> SET max_parallel_workers_per_gather =0;
SET
```

```


```

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_big1 WHERE
name='l_test';
QUERY PLAN
-----
Seq Scan on test_big1 (cost=0.00..991664.00 rows=1
width=25)
    (actual time=0.022..5329.100 rows=1 loops=1)
    Filter: ((name)::text = 'l_test'::text)
    Rows Removed by Filter: 49999999
    Planning time: 0.163 ms
    Execution time: 5329.136 ms
(5 rows)
```

```
SQL5329
3
```

6.2.2

index scan
Parallel index scan
test_big1

```
mydb=> EXPLAIN SELECT * FROM test_1 WHERE id=1;
              QUERY PLAN
```

```
-----
Index Scan using test_1_pkey on test_1 (cost=0.43..4.45
rows=1 width=26)
  Index Cond: (id = 1)
(2 rows)
```

Index Scan using
Parallel index scan
test_big1

```
mydb=> CREATE INDEX idx_test_big1_id ON test_big1 USING btree
(id);
CREATE INDEX
```

SQL ID 1

```
mydb=> EXPLAIN ANALYZE SELECT count(name) FROM test_big1
WHERE id<10000000;
```

QUERY PLAN

```
-----
Finalize Aggregate (cost=236183.98..236183.99 rows=1
width=8)
  (actual time=753.392..753.392 rows=1 loops=1)
  -> Gather (cost=236183.96..236183.97 rows=4 width=8)
      (actual time=750.133..753.384 rows=5 loops=1)
        Workers Planned: 4
        Workers Launched: 4
```


2636
2636



PostgreSQL10
SQL
CPU
IO
PostgreSQL 10
btree

6.2.3 index-only

index-only
index-only
SQL
ID
100

```
mydb=> SET max_parallel_workers_per_gather =0;  
SET
```

SQL

```
mydb=> EXPLAIN SELECT count(*) FROM test_big1 WHERE  
id<1000000;  
QUERY PLAN
```

```
-----  
Aggregate  (cost=36060.91..36060.92 rows=1 width=8)  
->  Index Only Scan using idx_test_big1_id on
```



```
test_big1 (cost=0.56..33313.99 rows=1098767 width=0)
      Index Cond: (id < 1000000)
(3 rows)
```

Index Only Scan
index-only
EXPLAIN ANALYZE SQL

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1 WHERE
id<1000000;
```

QUERY PLAN

```
-----
Aggregate (cost=35969.89..35969.90 rows=1 width=8)
  (actual time=253.571..253.571 rows=1 loops=1)
    -> Index Only Scan using idx_test_big1_id on
test_big1 (cost=0.56..33232.22 rows=1095066 width=0) (actual
time=0.038..0.179 rows=999999 loops=1)
      Index Cond: (id < 1000000)
      Heap Fetches: 999999
Planning time: 0.103 ms
Execution time: 253.617 ms
(6 rows)
```

253 index-only
index-only index-only

```
mydb=> SET max_parallel_workers_per_gather TO default;
SET
```

Parallel Index Only Scan

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1 WHERE  
id<1000000;
```

QUERY PLAN

```
-----  
-----  
      Finalize Aggregate  (cost=26703.66..26703.67 rows=1  
width=8)  
    (actual time=81.285..81.285 rows=1 loops=1)  
    -> Gather  (cost=26703.64..26703.65 rows=4 width=8)  
        (actual time=81.121..81.277 rows=5 loops=1)  
        Workers Planned: 4  
        Workers Launched: 4  
        -> Partial Aggregate  (cost=25703.64..25703.65  
rows=1 width=8) (actual time=75.778..75.778 rows=1 loops=5)  
            -> Parallel Index Only Scan using  
idx_test_big1_id on test_big1  (cost=0.56..25019.22  
rows=273766 width=0)  
                (actual time=0.045..59.398 rows=200000  
loops=5)  
                Index Cond: (id < 1000000)  
                Heap Fetches: 183366  
      Planning time: 0.113 ms  
      Execution time: 83.364 ms  
(10 rows)
```

Parallel Index Only Scan
index-only 83

6.2.4 bitmap heap

bitmap heap
Index Bitmap Heap SQL where

or Bitmap Index

```
mydb=> EXPLAIN SELECT * FROM test_big1 WHERE id=1 OR id=2;
```

QUERY PLAN

```
-----  
-----  
Bitmap Heap Scan on test_big1 (cost=5.15..9.17 rows=2  
width=25)
```

```
  Recheck Cond: ((id = 1) OR (id = 2))
```

```
    -> BitmapOr (cost=5.15..5.15 rows=2 width=0)
```

```
      -> Bitmap Index Scan on idx_test_big1_id  
(cost=0.00..2.57 rows=1 width=0)
```

```
        Index Cond: (id = 1)
```

```
      -> Bitmap Index Scan on idx_test_big1_id  
(cost=0.00..2.57 rows=1 width=0)
```

```
        Index Cond: (id = 2)
```

```
(7 rows)
```

Bitmap Index
Bitmap Index
test_big1 Bitmap Heap Bitmap Heap
SQL ID

```
EXPLAIN ANALYZE SELECT count(*) FROM test_big1 WHERE id  
<1000000 OR id > 49000000;
```

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1 WHERE  
id <1000000 OR id > 49000000;
```

QUERY PLAN

```
-----  
-----  
Finalize Aggregate (cost=406220.88..406220.89 rows=1
```

```

width=8)
  (actual time=241.186..241.186 rows=1 loops=1)
  -> Gather (cost=406220.46..406220.87 rows=4 width=8)
      (actual time=241.033..241.174 rows=5 loops=1)
      Workers Planned: 4
      Workers Launched: 4
      -> Partial Aggregate (cost=405220.46..405220.47
rows=1 width=8) (actual time=237.266..237.266 rows=1 loops=5)
          -> Parallel Bitmap Heap Scan on test_big1
(cost=28053.14..403933.27 rows=514876 width=0) (actual
time=83.059..189.073 rows=400000 loops=5)
              Recheck Cond: ((id < 1000000) OR (id >
49000000))
                  Heap Blocks: exact=2982
                  -> BitmapOr (cost=28053.14..28053.14
rows=2081101 width=0) (actual time=83.657..83.657 rows=0
loops=1)
                      -> Bitmap Index Scan on idx_test_big1_id
(cost=0.00..14219.03 rows=1094996 width=0)
                          (actual time=42.187..42.187
rows=999999 loops=1)
                              Index Cond: (id < 1000000)
                              -> Bitmap Index Scan on idx_test_big1_id
(cost=0.00..12804.35 rows=986105
width=0)
                                  (actual time=41.467..41.467
rows=1000000 loops=1)
                                      Index Cond: (id > 49000000)
Planning time: 0.157 ms
Execution time: 242.824 ms
(15 rows)

```

```

          Bitmap Heap
424ms

```

```

mydb=> SET max_parallel_workers_per_gather =0;
SET

```

```

SQL

```

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1 WHERE
id <1000000 OR id > 49000000;
```

QUERY PLAN

```
-----
Aggregate  (cost=432494.42..432494.43 rows=1 width=8) (actual
time=466.273.. 466.273 rows=1 loops=1)
  -> Bitmap Heap Scan on test_big1
      (cost=28053.14..427345.66 rows=2059505 width=0) (actual
time=86.859..299.822 rows=1999999 loops=1)
        Recheck Cond: ((id < 1000000) OR (id > 49000000))
        Heap Blocks: exact=13724
        -> BitmapOr  (cost=28053.14..28053.14 rows=2081101
width=0)
            (actual time=84.782..84.782 rows=0 loops=1)
            -> Bitmap Index Scan on idx_test_big1_id
                (cost=0.00..14219.03 rows= 1094996 width=0) (actual
time=42.704..42.704 rows=999999 loops=1)
                    Index Cond: (id < 1000000)
            -> Bitmap Index Scan on idx_test_big1_id
                (cost=0.00..12804.35 rows= 986105 width=0) (actual
time=42.076..42.076 rows=1000000 loops=1)
                    Index Cond: (id > 49000000)
        Planning time: 0.152 ms
        Execution time: 466.323 ms
(11 rows)
```

Bitmap Heap
466

6.3 性能调优

在PostgreSQL中，使用count和sum函数时，SQL引擎会生成不同的执行计划。本文将通过EXPLAIN ANALYZE来比较这两种聚合函数的性能差异。

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1;
QUERY PLAN
-----
Finalize Aggregate  (cost=525335.91..525335.92 rows=1
width=8)
  (actual time=2468.593.. 2468.594 rows=1 loops=1)
    -> Gather  (cost=525335.89..525335.90 rows=4 width=8)
          (actual time=2468. 386..2468.585 rows=5 loops=1)
            Workers Planned: 4
            Workers Launched: 4
            -> Partial Aggregate  (cost=524335.89..524335.90
rows=1 width=8)
                  (actual time=2463.532..2463.532 rows=1
loops=5)
                  -> Parallel Seq Scan on test_big1
(cost=0.00..493083.91
                      rows= 12500791 width=0) (actual
time=0.019..1456.506 rows=10000000
                      loops=5)
Planning time: 0.089 ms
Execution time: 2474.970 ms
(8 rows)
```

从执行计划中可以看出，使用count(*)时，SQL引擎会生成一个Partial Aggregate计划，然后是一个Finalize Aggregate计划。而使用sum(*)时，SQL引擎会生成一个Parallel Seq Scan计划。通过EXPLAIN ANALYZE，我们可以清楚地看到这两种聚合函数的性能差异。在6-1中，我们进一步探讨了如何优化SQL查询的性能。

count

```
mydb=> SET max_parallel_workers_per_gather = 0 ;
SET
```

```
top - 22:06:57 up 177 days,  8:01,  3 users,  load average: 1.98, 0.60, 0.21
Tasks: 213 total,   6 running, 207 sleeping,   0 stopped,   0 zombie
Cpu(s): 88.3%us,  7.4%sy,  0.0%ni,  4.3%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Mem:   8062340k total, 6643848k used, 1418492k free,   76912k buffers
Swap:      0k total,    0k used,    0k free, 5949344k cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
28940	postgres	20	0	1233m	168m	167m	R	84.8	2.1	0:02.55	postgres: bgworker: parallel worker for PID 24988
28939	postgres	20	0	1233m	217m	216m	R	81.4	2.8	0:02.45	postgres: bgworker: parallel worker for PID 24988
24988	postgres	20	0	1246m	1.0g	987m	R	77.8	12.7	3:06.46	postgres: pguser mydb [local] EXPLAIN
28938	postgres	20	0	1233m	195m	194m	R	69.8	2.5	0:02.10	postgres: bgworker: parallel worker for PID 24988
28941	postgres	20	0	1233m	141m	140m	R	67.1	1.8	0:02.02	postgres: bgworker: parallel worker for PID 24988

6-1 top

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1;
QUERY PLAN
```

```
-----
Aggregate  (cost=993115.55..993115.56 rows=1 width=8)
  (actual time=8655.614.. 8655.615 rows=1 loops=1)
    -> Seq Scan on test_big1  (cost=0.00..868107.64
      rows=50003164 width=0) (actual time=0.019..4898.227
      rows=50000000 loops=1)
      Planning time: 0.106 ms
      Execution time: 8655.666 ms
(4 rows)
```

247488653

2. 실행 계획 (Execution Plan)

```
mydb=> SET max_parallel_workers_per_gather =2;
SET
```

실행 계획 (Execution Plan)

```
mydb=> EXPLAIN ANALYZE SELECT count(*) FROM test_big1;
               QUERY PLAN
-----
Finalize Aggregate  (cost=629509.16..629509.17 rows=1
width=8)
  (actual time=3305. 585..3305.585 rows=1 loops=1)
    -> Gather  (cost=629509.15..629509.16 rows=2 width=8)
          (actual time=3305. 512..3305.578 rows=3 loops=1)
            Workers Planned: 2
            Workers Launched: 2
            -> Partial Aggregate  (cost=628509.15..628509.16
rows=1 width=8) (actual time=3302.362..3302.362 rows=1
loops=3)
              -> Parallel Seq Scan on test_big1
(cost=0.00..576422.52 rows= 20834652 width=0) (actual
time=0.021..2000.427 rows=16666667 loops=3)
                Planning time: 0.112 ms
                Execution time: 3314.371 ms
                (8 rows)
```

실행 계획 (Execution Plan) 결과: 3314ms, 68 SQL, 6-1

6-1 실행 계획 (Execution Plan)

$$\min_{\mathbf{w}} \max_{\mathbf{z}} \mathcal{L}(\mathbf{w}, \mathbf{z})$$

6.4 □□□□

```

    nested loop
join hash join

```

6.4.1 Nested loop

□□□□Nested loop□□□□□□□□□□□□□□□□□□
□□

```
for (i = 0; i < length(outer); i++)
    for (j = 0; j < length(inner); j++)
        if (outer[i] == inner[j])
            output(outer[i], inner[j]);
```

```

Nested loop
test_small

```

```
CREATE TABLE test_small(id int4, name character
varying(32));

INSERT INTO test_small(id,name)
SELECT n, n|| ' _small' FROM generate_series(1,8000000) n ;
```

□□□□□□□□□□□□□□□□

```
mydb=> CREATE INDEX idx_test_small_id ON test_small USING
btree (id);
CREATE INDEX

mydb=> ANALYZE test_small;
ANALYZE
```

ANALYZE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□

```
mydb=> EXPLAIN ANALYZE SELECT test_small.name
      FROM test_big1, test_small
      WHERE test_big1.id = test_small.id
      AND test_small.id < 10000;
QUERY PLAN
-----
-----
Gather  (cost=1138.18..31628.76 rows=10217 width=13) (actual
time=1.036..16.207 rows=9999 loops=1)
  Workers Planned: 3
  Workers Launched: 3
  -> Nested Loop  (cost=138.18..29607.06 rows=3296
width=13)
    (actual time= 0.244..10.895 rows=2500 loops=4)
    -> Parallel Bitmap Heap Scan on test_small
(cost=137.61..14796.65
      rows= 3296 width=17) (actual time=0.203..0.653
rows=2500 loops=4)
      Recheck Cond: (id < 10000)
      Heap Blocks: exact=10
      -> Bitmap Index Scan on idx_test_small_id
(cost=0.00..135.06 rows=10217 width=0) (actual
time=0.652..0.652 rows=9999 loops=1)
        Index Cond: (id < 10000)
        -> Index Only Scan using idx_test_big1_id on
test_big1  (cost=0.56..4.48 rows=1 width=4) (actual
```

```
time=0.003..0.003 rows=1 loops=9999)
      Index Cond: (id = test_small.id)
      Heap Fetches: 1674
Planning time: 0.427 ms
Execution time: 17.548 ms
(14 rows)
```

test_big1 Index
Only id 10000 Nested
loop test_small Bitmap
Heap id 10000
SQL 17.5

```
mydb=> SET max_parallel_workers_per_gather =0;
SET
```

```
mydb=> EXPLAIN ANALYZE SELECT test_small.name
      FROM test_big1, test_small
      WHERE test_big1.id = test_small.id
      AND test_small.id < 10000;
      QUERY PLAN
```

```
-----
Nested Loop (cost=1.00..46213.80 rows=10217 width=13)
(actual time=0.025..29.054 rows=9999 loops=1)
  -> Index Scan using idx_test_small_id on test_small
      (cost=0.43..304.23 rows=10217 width=17) (actual
time=0.014..2.284 rows=9999 loops=1)
      Index Cond: (id < 10000)
  -> Index Only Scan using idx_test_big1_id on test_big1
      (cost=0.56..4.48 rows=1 width=4) (actual time=0.002..0.002
rows=1 loops=9999)
```

```
Index Cond: (id = test_small.id)
Heap Fetches: 9999
Planning time: 0.262 ms
Execution time: 29.606 ms
(8 rows)
```

SQL29.606ms
8 rows

6.4.2 Merge join

Merge join
Merge join

```
mydb=> EXPLAIN ANALYZE SELECT test_small.name
FROM test_big1, test_small
WHERE test_big1.id = test_small.id
AND test_small.id < 200000;
```

QUERY PLAN

```
-----
Gather  (cost=1001.79..195146.04 rows=204565 width=13)
  (actual time=0.308.. 160.192 rows=199999 loops=1)
    Workers Planned: 4
    Workers Launched: 4
    -> Merge Join  (cost=1.79..173689.54 rows=51141
width=13)
      (actual time=4.059.. 127.164 rows=40000 loops=5)
        Merge Cond: (test_big1.id = test_small.id)
        -> Parallel Index Only Scan using
idx_test_big1_id on test_big1  (cost=0.56..1017275.56
rows=12500000 width=4)
          (actual time=0.044..16.223 rows=40001
loops=5)
          Heap Fetches: 50875
        -> Index Scan using idx_test_small_id on
```

```
test_small  (cost=0.43..6010.32
              rows=204565 width=17) (actual
time=0.034..64.427 rows=199999 loops=5)
              Index Cond: (id < 200000)
  Planning time: 0.255 ms
  Execution time: 171.755 ms
(11 rows)
```

```

      171
11
```

```
mydb=> SET max_parallel_workers_per_gather =0;
SET
```

```


```

```
mydb=> EXPLAIN ANALYZE SELECT test_small.name
      FROM test_big1, test_small
      WHERE test_big1.id = test_small.id
      AND test_small.id < 200000;
              QUERY PLAN
```

```
-----
Merge Join  (cost=1.79..249728.34 rows=204565 width=13)
(actual time=0.034..198.331 rows=199999 loops=1)
  Merge Cond: (test_big1.id = test_small.id)
    -> Index Only Scan using idx_test_big1_id on test_big1
        (cost=0.56..1392275.56 rows=50000000 width=4) (actual
time=0.018..60.445
        rows=200000 loops=1)
        Heap Fetches: 200000
    -> Index Scan using idx_test_small_id on test_small
        (cost=0.43..6010.32 rows=204565 width=17) (actual
time=0.013..50.531 rows=199999 loops=1)
        Index Cond: (id < 200000)
  Planning time: 0.264 ms
```

Execution time: 209.387 ms
(8 rows)

PostgreSQL 9.5.4 SQL Server 2008 R2
Query Plan

6.4.3 Hash join

PostgreSQL 9.5.4 Hash join
Query Plan
Hash join
Hash join
Hash join

```
mydb=> DROP INDEX idx_test_big1_id;  
DROP INDEX  
mydb=> DROP INDEX idx_test_small_id ;  
DROP INDEX  
  
mydb=> SET max_parallel_workers_per_gather =0;  
SET
```

Query Plan

```
mydb=> EXPLAIN SELECT test_small.name  
        FROM test_big1 JOIN test_small ON (test_big1.id =  
test_small.id)  
        AND test_small.id < 100;  
                                QUERY PLAN  
-----  
Hash Join  (cost=150871.78..1205043.77 rows=800 width=13)  
  Hash Cond: (test_big1.id = test_small.id)  
    -> Seq Scan on test_big1  (cost=0.00..866664.00
```



```
rows=50000000 width=4)
  -> Hash (cost=150861.78..150861.78 rows=800 width=17)
      -> Seq Scan on test_small
(cost=0.00..150861.78 rows=800 width=17)
      Filter: (id < 100)
(6 rows)
```

SQL

```
mydb=> EXPLAIN ANALYZE SELECT test_small.name
      FROM test_big1 JOIN test_small ON (test_big1.id =
test_small.id)
      AND test_small.id < 100;
```

QUERY PLAN

```
-----
Hash Join (cost=151317.06..1206944.59 rows=802 width=13)
  (actual time=735.778..11259.744 rows=99 loops=1)
  Hash Cond: (test_big1.id = test_small.id)
    -> Seq Scan on test_big1 (cost=0.00..868107.64
rows=50003164 width=4) (actual time=0.015..5424.515
rows=50000000 loops=1)
    -> Hash (cost=151307.04..151307.04 rows=802 width=17)
        (actual time=735.750..735.750 rows=99 loops=1)
        Buckets: 1024 Batches: 1 Memory Usage: 13kB
        -> Seq Scan on test_small
(cost=0.00..151307.04 rows=802 width=17)
            (actual time=0.010..735.731 rows=99 loops=1)
            Filter: (id < 100)
            Rows Removed by Filter: 7999901
Planning time: 0.134 ms
Execution time: 11259.789 ms
(10 rows)
```

11259 4

```
mydb=> SET max_parallel_workers_per_gather =4;
SET
```

SQL

```
mydb=> EXPLAIN ANALYZE SELECT test_small.name
      FROM test_big1 JOIN test_small ON (test_big1.id =
test_small.id)
      AND test_small.id < 100;
```

QUERY PLAN

```
-----
-----
Gather  (cost=152317.06..692280.94 rows=802 width=13)
  (actual time=1152.263..4304.757 rows=99 loops=1)
    Workers Planned: 4
    Workers Launched: 4
    -> Hash Join  (cost=151317.06..691280.94 rows=200
width=13) (actual time= 3667.973..4298.423 rows=20 loops=5)
      Hash Cond: (test_big1.id = test_small.id)
      -> Parallel Seq Scan on test_big1
(cost=0.00..493083.91 rows=12500791 width=4) (actual
time=0.025..1737.558 rows=10000000 loops=5)
      -> Hash  (cost=151307.04..151307.04 rows=802
width=17)
        (actual time=1106.665..1106.665 rows=99
loops=5)
        Buckets: 1024  Batches: 1  Memory Usage:
13kB
        -> Seq Scan on test_small
(cost=0.00..151307.04 rows=802 width=17)
          (actual time=863.670..1106.624 rows=99
loops=5)
          Filter: (id < 100)
          Rows Removed by Filter: 7999901
Planning time: 0.156 ms
Execution time: 4315.928 ms
(13 rows)
```

4315

6.5 索引

PostgreSQL 索引类型包括 B-tree、hash、index-only、bitmap heap 等。索引是数据库中最常用的优化手段之一，可以显著提高查询效率。索引的创建和维护成本较高，因此需要根据实际情况合理选择索引类型和索引策略。索引的优化是数据库性能优化的重要组成部分，也是数据库管理员必须掌握的技能之一。

07 00000000

Lock-Based Concurrency Control

MVCC

Multi-Version Concurrency Control

PostgreSQL

PostgreSQL MVCC

7.1 □□□□□□□□

7.1.1 〇〇〇〇〇〇〇〇〇〇〇〇

```

--
-- PostgreSQL
BEGIN...END/COMMIT/ROLLBACK

```

- Atomicity

· Consistency

· **Isolation**

Isolation is a technique used to prevent a transaction from seeing uncommitted changes made by other transactions. It is achieved by locking the data being accessed, ensuring that only one transaction can modify it at a time. This prevents dirty reads, non-repeatable reads, and phantom reads.

Isolation levels define the degree of isolation, ranging from Read Uncommitted (lowest) to Serializable (highest). The levels are:

- Read Uncommitted: Allows dirty reads.
- Read Committed: Prevents dirty reads.
- Repeatable Read: Prevents dirty reads and non-repeatable reads.
- Serializable: Prevents dirty reads, non-repeatable reads, and phantom reads.

Isolation is implemented using various locking mechanisms, such as two-phase locking (2PL) and multiversion concurrency control (MVCC).

· **Durability** 데이터가 영구적으로 저장되는 것을 보장하는 것

ACID
WAL
MVCC

7.1.2 데이터 일관성

데이터 일관성은 트랜잭션이 실행되는 동안 데이터베이스의 상태를 유지하는 것을 보장하는 것
PostgreSQL
Dirty read
Non-repeatable read
Phantom Read
Serialization Anomaly



PostgreSQL
READ
UNCOMMITTED
READ COMMITTED
PostgreSQL
READ
COMMITTED

PostgreSQLのREAD COMMITTED の挙動

1. 概要

PostgreSQLのREAD COMMITTEDは、
INSERT、UPDATE、DELETEの書き込みは
ROLLBACKしても書き込みは成功する。
読み取りは、書き込みが完了した
後にのみ行われる。

以下は、PostgreSQLのREAD COMMITTEDの挙動を示す。

```
CREATE TABLE tbl_mvcc (  
  id INT NOT NULL AUTO_INCREMENT,  
  ival INT,  
  PRIMARY KEY (id)  
);  
-- 書き込み  
INSERT INTO tbl_mvcc (ival) VALUES (1);
```

以下は、T1とT2の同時実行を示す。

T1	T2
<pre>set session transaction isolation level read uncommitted;start transaction; select * from tbl_mvcc where id = 1; /* id=1,ival=1 */</pre>	

(续)

T1	T2
	start transaction;
	update tbl_mvcc set ival = 10 where id = 1;
select * from tbl_mvcc where id = 1; /* id=1,ival=10 */	
	rollback;

```

--T1tbl_mvcc--
id=1ival=1--T2--id=1--ival
10--T1tbl_mvcc--T2--ival
--1--T1--ival--10--T2--
--ROLLBACK--T1--

```

2. 問題

[illegible]

```
CREATE TABLE tbl_mvcc (
    id SERIAL PRIMARY KEY,
    ival INT
);
-- [][][][]
INSERT INTO tbl_mvcc (ival) VALUES (1);
```


T1 T2

T1	T2
BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED;	
SELECT id,ival FROM tbl_mvcc WHERE id = 1;	
<pre> id ival ---+----- 1 1 (1 row) </pre>	
	<pre> BEGIN; UPDATE tbl_mvcc SET ival = 10 WHERE id = 1; COMMIT; </pre>
<pre> SELECT id,ival FROM tbl_mvcc WHERE id = 1; id ival ---+----- 1 10 (1 row) </pre>	
END;	

```

-- T1 tbl_mvcc id=1 ival=1
-- T2 id=1 ival=10
-- T2 COMMIT
-- T1 tbl_mvcc ival=10
-- SELECT T2
-- T2 T1

```

3.

INSERT DELETE

tbl_mvcc

T1	T2
<pre>BEGIN TRANSACTION ISOLATION LEVEL READ COMMITTED; SELECT id,ival FROM tbl_mvcc WHERE id > 3 AND id < 10; id ival ---+----- 4 4 5 5 (2 rows)</pre>	
	<pre>BEGIN; INSERT INTO tbl_mvcc (id , ival) VALUES (6 , 6); END;</pre>
<pre>SELECT id,ival FROM tbl_mvcc WHERE id > 3 AND id < 10; id ival ---+----- 4 4 5 5 6 6 (3 rows)</pre>	
<pre>END;</pre>	

T1tbl_mvccid3
10T2id6
T1WHEREid3
10T1
UPDATEINSERT
DELETE

7.1.3 ANSI SQL隔离级别

数据库的隔离级别是指数据库在并发操作时，对事务的隔离程度。ANSI SQL定义了四种隔离级别，从低到高分别是：Read Uncommitted、Read Committed、Repeatable Read和Serializable。这四种隔离级别分别对应不同的并发控制策略，如锁、多版本控制等。不同的数据库系统对ANSI SQL隔离级别的支持程度不同，有些数据库可能只支持其中的一部分。

·Read Uncommitted 读未提交。这是最低的隔离级别，允许事务读取尚未提交的数据。在这种级别下，事务可以读到其他事务尚未提交的修改，这可能导致脏读。大多数数据库系统都不支持这个级别，因为它会破坏数据的一致性。

·Read Committed 读已提交。这是默认的隔离级别，保证事务只能看到已经提交的数据。在这种级别下，事务不会读到其他事务尚未提交的修改，这可以避免脏读。但是，它仍然可能遇到不可重复读的问题。

·Repeatable Read 可重复读。这种隔离级别保证事务在多次读取同一数据时，得到的结果是一致的。在这种级别下，事务不会遇到不可重复读的问题，但仍然可能遇到幻读的问题。

·Serializable 可串行化。这是最高的隔离级别，保证事务的执行结果与串行执行的结果一致。在这种级别下，事务不会遇到任何并发问题，但它的性能通常是最差的，因为它需要大量的锁来保证串行化。

ANSI SQL

隔离级别	脏 读	不可重复读	幻 读
Read Uncommitted	可能	可能	可能
Read Committed	不可能	可能	可能
Repeatable Read	不可能	不可能	可能
Serializable	不可能	不可能	不可能

Read Committed

7.2 PostgreSQL隔离级别

数据库隔离级别是数据库系统实现事务并发控制的重要机制。SQL标准定义了四种隔离级别：“Read Uncommitted”、“Read Committed”、“Repeatable Read”和“Serializable”。PostgreSQL支持这四种隔离级别，并默认使用“Read Committed”级别。PostgreSQL的“Read Uncommitted”级别允许脏读，即一个事务可以读取另一个未提交事务的数据。PostgreSQL的“Read Committed”级别防止脏读，但允许不可重复读。PostgreSQL的“Repeatable Read”级别防止脏读和不可重复读，但允许幻读。PostgreSQL的“Serializable”级别防止脏读、不可重复读和幻读，保证事务的串行化执行。T1和T2是两个事务。

Repeatably Readable Transaction

T1	T2
<pre>BEGIN TRANSACTION ISOLATION LEVEL REPEATABLE READ; SELECT id,ival FROM tbl_mvcc WHERE id > 3 AND id < 10; id ival ---+----- 4 4 5 5 (2 rows)</pre>	
	<pre>BEGIN; INSERT INTO tbl_mvcc (id , ival) VALUES (6 , 6); END;</pre>
<pre>SELECT id,ival FROM tbl_mvcc WHERE id > 3 AND id < 10; id ival ---+----- 4 4 5 5 (2 rows)</pre>	
END;	

Transaction T1 reads rows with id between 3 and 10 from table tbl_mvcc. Transaction T2 inserts a new row with id 6 and ival 6. Transaction T1 reads the same rows again. The results are the same as the first read. This is because the transaction is repeatable read.

7.2.1 Serializable

PostgreSQL Serializable Read Committed

```
mydb=# SELECT name,setting FROM pg_settings WHERE name =  
'default_transaction_isolation';
```

```
          name          |      setting        
-----+-----  
 default_transaction_isolation | repeatable read  
(1 row)
```

```
mydb=# SELECT  
current_setting('default_transaction_isolation');  
      current_setting
```

```
-----  
 repeatable read  
(1 row)
```

7.2.2 配置数据库隔离级别

1 编辑 postgresql.conf 文件
default_transaction_isolation 配置为 repeatable read
并执行 reload 命令

2 执行 ALTER SYSTEM 命令

```
mydb=# ALTER SYSTEM SET default_transaction_isolation TO  
'REPEATABLE READ';
```

```
ALTER SYSTEM
```

```
mydb=# SELECT pg_reload_conf();  
      pg_reload_conf
```

```
-----  
 t
```

```
(1 row)
```

```
mydb=# SELECT current_setting('transaction_isolation');  
      current_setting
```

```
-----  
 repeatable read  
(1 row)
```

7.2.3 数据库事务隔离级别

数据库事务隔离级别

```
mydb=# SHOW transaction_isolation ;
      transaction_isolation
-----
      read committed
(1 row)
```

□

```
mydb=# SELECT current_setting('transaction_isolation');
      current_setting
-----
      read committed
(1 row)
```

7.2.4 数据库事务隔离级别

数据库事务隔离级别

```
mydb=# SET SESSION CHARACTERISTICS AS TRANSACTION ISOLATION
      LEVEL READ UNCOMMITTED;
SET
mydb=# SHOW transaction_isolation ;
      transaction_isolation
-----
      read uncommitted
(1 row)
```

7.2.5 数据库事务隔离

数据库事务隔离

```
mydb=# START TRANSACTION ISOLATION LEVEL READ UNCOMMITTED;
START TRANSACTION
...
...
...
mydb=# END;
COMMIT
```

数据库

```
mydb=# BEGIN ISOLATION LEVEL READ UNCOMMITTED READ WRITE;
...
...
...
mydb=# END/COMMIT/ROLLBACK;
```

START TRANSACTION BEGIN 数据库事务
数据库事务

7.3 PostgreSQL

PostgreSQL implements two different concurrency control mechanisms: Lock-Based Concurrency Control (LBCC) and Multi-Version Concurrency Control (MVCC). LBCC is a pessimistic concurrency control mechanism that uses locks to prevent conflicts between transactions. MVCC is an optimistic concurrency control mechanism that uses versioning to allow multiple versions of a row to exist simultaneously. PostgreSQL also supports Optimistic Concurrency Control (OCC) and Pessimistic Concurrency Control (PCC).

7.3.1 Locks

PostgreSQL implements two types of locks: Exclusive locks (X) and Share locks (S).

数据库事务的隔离性是指一个事务的执行过程不受其他事务的干扰，即一个事务在执行过程中，其数据不会被其他事务所修改。

数据库事务的持久性是指一个事务一旦提交，其数据就会永久地存储在数据库中，即使系统发生故障，数据也不会丢失。

数据库事务的粒度（Granularity）是指事务中包含的操作的数量。粒度越小，事务的原子性越强，但性能越低。粒度越大，事务的性能越高，但原子性越弱。数据库事务的锁定协议（Locking Protocol）是指数据库系统为了保证事务的隔离性而采用的一种机制。常见的锁定协议有悲观锁（Pessimistic Locking）和乐观锁（Optimistic Locking）。

PostgreSQL 数据库支持 Advisory Lock（ Advisory Lock 是一种用户自定义的锁，用于实现分布式事务的协调。

7.3.2 数据库事务的隔离性

数据库事务的隔离性是指一个事务的执行过程不受其他事务的干扰。数据库系统通过多种机制来实现事务的隔离性，其中最常用的是 MVCC（Multi-Version Concurrency Control）。

MVCC
Oracle
PostgreSQL MySQL Innodb MVCC
MVCC
MVCC

MVCC
READ COMMITTED
SQL
SERIALIZABLE

PostgreSQL int32
IDxid
idid
PostgreSQL
PostgreSQL
PostgreSQL
4
xminxmaxcmincmaxcmincmax
xminxmax

PostgreSQLのMVCCについて
PostgreSQLのMVCCについて

```
mydb=# SELECT xmin,xmax,cmin,cmax,id,ival FROM tbl_mvcc
WHERE id = 1;
```

xmin	xmax	cmin	cmax	id	ival
1930	0	0	0	1	1

(1 row)

xminは、このレコードが作成されたときに、
xidはPostgreSQLの内部で生成されたxminはxmax
は、このレコードが作成されたときに、

1. xminについて

PostgreSQLの内部で生成されたxidはxminはxmax
は、このレコードが作成されたときに、xminはxmax
は、このレコードが作成されたときに、

1. xminについて
xminは、このレコードが作成されたときに、

```
mydb=# BEGIN;
mydb=# SELECT txid_current();
txid_current
-----
1937
(1 row)
mydb=# INSERT INTO tbl_mvcc(id,ival) VALUES(7,7);
INSERT 0 1
```

```
mydb=# SELECT xmin,xmax,cmin,cmax,id,ival FROM tbl_mvcc
WHERE id = 7;
```

xmin	xmax	cmin	cmax	id	ival
1937	0	0	0	7	7

(1 row)

```
SELECT txid_current()
xid 1937 id 7
xmin 1937 xid
```

```
mydb=# BEGIN;
BEGIN
mydb=# SELECT txid_current();
txid_current
```

1938

(1 row)

```
mydb=# SELECT * FROM tbl_mvcc WHERE id = 7;
id | ival
----+-----
(0 rows)
```

```
mydb=# END;
COMMIT
```

2. Repeating the above steps with the isolation level set to Repeatable Read and Serializable.

mysql> SET TRANSACTION ISOLATION LEVEL REPEATABLE READ;

mysql> BEGIN;

mysql> SELECT txid_current();

txid_current

1939

(1 row)

mysql> SET TRANSACTION ISOLATION LEVEL SERIALIZABLE;

mysql> BEGIN;

mysql> SELECT txid_current();

txid_current

1939

(1 row)

mysql> BEGIN;

mysql> BEGIN;

mysql> SELECT txid_current();

txid_current

1940

(1 row)

mysql> INSERT INTO tbl_mvcc (id,ival) VALUES (7,7);

INSERT 0 1

mysql> SELECT xmin,xmax,cmin,cmax,id,ival FROM tbl_mvcc

WHERE id = 7;

xmin	xmax	cmin	cmax	id	ival
1940	0	0	0	7	7

(1 row)

mysql> COMMIT;

COMMIT

xxxxxxxid1940xxxxxxxxxxxxxxxx
xxxminxxxxxxxxxxidxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxx

```
mydb=# SELECT xmin,xmax,cmin,cmax,id,ival FROM tbl_mvcc
WHERE id = 7;
   xmin | xmax | cmin | cmax | id | ival
-----+-----+-----+-----+---+----
(0 rows)
mydb=# END;
COMMIT
```

xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxXID1939xx
xxxxxxxxxxxmin1940xxxxxxxxxxxminxx
xxid7xxxxxxxxxxxxxxxx

2. xxxmaxxxxxxxxx

xxxmaxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxx1xxxxxxxxxmaxxxxxxxxx2
xxxxxxxxxxxxxxxxxidxxxxxxxx3xxxxxx
xxxxxxxxxxxxxxxxCOMMITROLLBACKxxxxxid
xxxxxxxxxxxxxxxx4xxxxxxxxxxxxxxxxxidxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx

7.3.3 xxpageinspectxxMVCC

PostgreSQL
 PostgreSQL pageinspect
 pageinspect
 pageinspect

```

mydb=# CREATE EXTENSION pageinspect;
CREATE EXTENSION
mydb=# \dx+ pageinspect
        Objects in extension "pageinspect"
        Object description
-----

```

```

...
function get_raw_page(text,integer)
...
function heap_page_items(bytea)
...
(19 rows)

```

get_raw_page
 get_raw_page relname text fork text
 blkno int
 get_raw_page relname
 text blkno int
 relation
 relname relation name fork main
 vm fsm init fork main main
 vm fsm free space
 map init
 get_raw_page
 bytea
 heap_page_items
 heap_page_items
 heap_page_items

MVCC
get_raw_page
heap_page_items

PostgreSQL MVCC
BRUCE MOMJIAN
<http://momjian.us/>

```
DROP VIEW IF EXISTS v_pageinspect;  
CREATE VIEW v_pageinspect AS  
SELECT  '(0,' || lp || ')' AS ctid,  
        CASE lp_flags  
            WHEN 0 THEN 'Unused'  
            WHEN 1 THEN 'Normal'  
            WHEN 2 THEN 'Redirect to ' || lp_off  
            WHEN 3 THEN 'Dead'  
        END,  
        t_xmin::text::int8 AS xmin,  
        t_xmax::text::int8 AS xmax,  
        t_ctid  
FROM heap_page_items(get_raw_page('tbl_mvcc', 0))  
ORDER BY lp;
```

INSERT
xminxidxmaxNULL

```
mydb=# BEGIN;  
mydb=# SELECT txid_current()  
txid_current  
-----  
565  
(1 row)
```

```
-- 1 id 565
mydb=# INSERT INTO tbl_mvcc (ival) VALUES (1);
mydb=# SELECT * FROM v_pageinspect;
   ctid | case | xmin | xmax | t_ctid
-----+-----+-----+-----+-----
   (0,1) | Normal | 565 | 0 | (0,1)
(1 row)
mydb=# END;
-- 2 INSERT xminxid565xmaxNULL
```

DELETE xmaxxid

```
mydb=# BEGIN;
BEGIN
mydb=# SELECT txid_current();
   txid_current
-----
          561
(1 row)

mydb=# DELETE FROM tbl_mvcc WHERE id = 1;
DELETE 1
mydb=# SELECT * FROM v_pageinspect;
   ctid | case | xmin | xmax | t_ctid
-----+-----+-----+-----+-----
   (0,1) | Normal | 565 | 566 | (0,1)
(1 row)

mydb=# END;
COMMIT
```

UPDATEDELETE

INSERT

```

mydb=# INSERT INTO tbl_mvcc (ival) VALUES (2); -- 
UPDATE
mydb=# BEGIN;
mydb=# SELECT txid_current();
        txid_current
-----
        639
(1 row)
-- id=639
mydb=# SELECT * FROM tbl_mvcc;
   id | ival
-----+-----
    2 |    2
(1 row)
-- tbl_mvcc
mydb=# SELECT * FROM v_pageinspect;
   ctid | case | xmin | xmax | t_ctid
-----+-----+-----+-----+-----
(0,1) | Normal | 567 |    0 | (0,1)
(1 row)
-- pageinspect page xmin xid
xid=567
mydb=# UPDATE tbl_mvcc SET ival = 20 WHERE id = 2;
-- 
mydb=# SELECT * FROM v_pageinspect;
   ctid | case | xmin | xmax | t_ctid
-----+-----+-----+-----+-----
(0,1) | Normal | 567 | 639 | (0,2)
(0,2) | Normal | 639 |    0 | (0,2)
(2 rows)
mydb=# END;

```

pageinspect page
 UPDATE DELETE INSERT
 xid=567 ctid=0
 1 xmin=567 xmax xid
 639 page ctid=0 2
 xmin xid=639

11

```
[root@pgghost1 ~]# yum install -y pg_repack10
```

pg_repack

```
mydb=# CREATE EXTENSION pg_repack;
```

```
pg_repack tbl_mvcc
```

```
[pos tgres@pghost1 ~]# /usr/pgsql-10/bin/pg_repack -t
tbl_mvcc -j 2 -D -k -h pghost1-U postgres -d mydb
```

PostgreSQL
 MVCC

7.3.5 数据库DDL

```
PostgreSQL
CREATE TABLE DDL
CREATE TABLE DDL
TRUNCATE
```

```

mydb=# DROP TABLE IF EXISTS tbl_test;
NOTICE:  table "tbl_test" does not exist, skipping
DROP TABLE
mydb=# BEGIN;
BEGIN
mydb=# CREATE TABLE tbl_test (ival int);
CREATE TABLE
mydb=# INSERT INTO tbl_test VALUES (1);
INSERT 0 1
mydb=# ROLLBACK;
ROLLBACK
mydb=# SELECT * FROM tbl_test;
ERROR:  relation "tbl_test" does not exist

```

TRUNCATE

```

mydb=# SELECT COUNT(*) FROM tbl_mvcc;
count
-----
9
(1 row)
mydb=# BEGIN;
BEGIN
mydb=# TRUNCATE tbl_mvcc ;
TRUNCATE TABLE
mydb=# ROLLBACK;
ROLLBACK
mydb=# SELECT COUNT(*) FROM tbl_mvcc;
count
-----
9
(1 row)

```

TRUNCATE

7.4 練習問題

[illegible]

8

Oracle
PostgreSQL
10 PostgreSQL
PostgreSQL10

8.1 数据库

数据库

- 数据库是存储在计算机中的有组织的数据集合

- 数据库操作包括 **DROP TABLE** 删除表和 **DELETE** 删除数据

- 数据库操作还包括插入、更新、查询和删除数据

数据库操作是数据库管理系统的核心功能

8.2 ☐☐☐☐☐

A diagram showing four rows of blocks representing a number sentence. Row 1: 10 blocks. Row 2: 10 blocks. Row 3: 10 blocks. Row 4: 10 blocks. The total number of blocks is 40.

8.2.1 練習問題

PostgreSQL

tbl_log

```
mydb=> CREATE TABLE tbl_log(id int4,create_date date,log_type
text);
CREATE TABLE
```

tbl_log_sql SQL

```
mydb=> CREATE TABLE tbl_log_sql(sql text) INHERITS(tbl_log);
CREATE TABLE
```

```

    INHERITS(tbl_log)tbl_log_sql
tbl_logsql

```

CREATE TABLE tbl_log_sql (

```
mydb=> \d tbl_log_sql
          Table "pguser.tbl_log_sql"
  Column | Type   | Collation | Nullable | Default
-----+-----+-----+-----+-----
    id   | integer |           |          |
 create_date | date   |           |          |
 log_type | text   |           |          |
    sql  | text   |           |          |
Inherits: tbl_log
```

CREATE TABLE tbl_log_sql (

tbl_log sql Inherits

tbl_log tbl_log

CREATE TABLE tbl_log (

tbl_log

```
mydb=> INSERT INTO tbl_log VALUES (1,'2017-08-26',null);
INSERT 0 1
mydb=> INSERT INTO tbl_log_sql VALUES(2,'2017-08-
27',null,'select 2');
INSERT 0 1
```

CREATE TABLE tbl_log (

```
mydb=> SELECT * FROM tbl_log;
 id | create_date | log_type
-----+-----+-----
  1 | 2017-08-26 |
```

```
2 | 2017-08-27 |
(2 rows)
```

SQLのOID

```
mydb=> SELECT tableoid, * FROM tbl_log;
      tableoid | id | create_date | log_type
-----+-----+-----+-----
      16854   |  1 | 2017-08-26   |
      16860   |  2 | 2017-08-27   |
(2 rows)
```

tableoid OID pg_class

```
mydb=> SELECT p.relname,c.*
      FROM tbl_log c, pg_class p
      WHERE c.tableoid = p.oid;
      relname | id | create_date | log_type
-----+-----+-----+-----
      tbl_log |  1 | 2017-08-26   |
      tbl_log_sql |  2 | 2017-08-27   |
(2 rows)
```

ONLY

```
mydb=> SELECT * FROM ONLY tbl_log;
      id | create_date | log_type
-----+-----+-----
```

1 | 2017-08-26 |
(1 row)

UPDATE DELETE SELECT
ONLY DML

```
mydb=> DELETE FROM tbl_log;
DELETE 2
mydb=> SELECT count(*) FROM tbl_log;
count
-----
0
(1 row)
```



UPDATE
DELETE DML

8.2.2

1


```

CREATE TABLE log_ins_history(CHECK ( create_time < '2017-01-01' )) INHERITS(log_ins);
CREATE TABLE log_ins_201701(CHECK ( create_time >= '2017-01-01' and create_time < '2017-02-01')) INHERITS(log_ins);
CREATE TABLE log_ins_201702(CHECK ( create_time >= '2017-02-01' and create_time < '2017-03-01')) INHERITS(log_ins);
...
CREATE TABLE log_ins_201712(CHECK ( create_time >= '2017-12-01' and create_time < '2018-01-01')) INHERITS(log_ins);

```

□□□□□□□□□□□□□□□□□□□□□□□□

```

CREATE INDEX idx_his_ctime ON log_ins_history USING btree
(create_time);
CREATE INDEX idx_log_ins_201701_ctime ON log_ins_201701 USING
btree (create_time);
CREATE INDEX idx_log_ins_201702_ctime ON log_ins_201702 USING
btree (create_time);
...
CREATE INDEX idx_log_ins_201712_ctime ON log_ins_201712 USING
btree (create_time);

```

□□□□□□□□□□□□□□□□□□□□□□□□

□□□

```

CREATE OR REPLACE FUNCTION log_ins_insert_trigger()
    RETURNS trigger
    LANGUAGE plpgsql
AS $function$
BEGIN
    IF ( NEW. create_time < '2017-01-01' ) THEN
        INSERT INTO log_ins_history VALUES (NEW.*);
    ELSIF ( NEW.create_time>='2017-01-01' and
NEW.create_time<'2017-02-01' ) THEN
        INSERT INTO log_ins_201701 VALUES (NEW.*);
    ELSIF ( NEW.create_time>='2017-02-01' and

```

```

NEW.create_time<'2017-03-01' ) THEN
    INSERT INTO log_ins_201702 VALUES (NEW.*);
    ...
ELSIF ( NEW.create_time>='2017-12-01' and
NEW.create_time<'2018-01-01' ) THEN
    INSERT INTO log_ins_201712 VALUES (NEW.*);
ELSE
    RAISE EXCEPTION 'create_time out of range. Fix the
log_ins_insert_trigger() function!';
END IF;
RETURN NULL;
END;
$function$;

```

new.*

```

CREATE TRIGGER insert_log_ins_trigger BEFORE INSERT ON
log_ins FOR EACH ROW EXECUTE PROCEDURE
log_ins_insert_trigger();

```

log_ins

log_ins_insert_trigger

DELETE UPDATE INSERT



user_id

user_id

8.2.3 数据插入

向log_ins表插入数据

```
INSERT INTO log_ins(user_id,create_time)
SELECT round(1000000000*random()),generate_series('2016-12-01'::date,
'2017-12-01'::date, '1 minute');
```

查看插入的数据

```
mydb=> SELECT * FROM log_ins LIMIT 2;
      id | user_id |      create_time
-----+-----+-----
  570242 | 24040985 | 2016-12-01 00:00:00
  570243 | 10814368 | 2016-12-01 00:01:00
(2 rows)
```

查看表的行数

```
mydb=> SELECT count(*) FROM ONLY log_ins;
 count
-----
      0
(1 row)
```

```
mydb=> SELECT count(*) FROM log_ins;
 count
-----
525601
(1 row)
```

[illegible]

```
mydb=> SELECT min(create_time),max(create_time) FROM
log_ins_201701;
      min                |                max
-----+-----
 2017-01-01 00:00:00 | 2017-01-31 23:59:00
(1 row)
```

[illegible]

```
mydb=> \dt+ log_ins*

```

List of relations					
Schema	Name	Type	Owner	Size	
Description					

```

-----+-----+-----+-----+-----+
pguser | log_ins          | table | pguser | 0 bytes      |
pguser | log_ins_201701   | table | pguser | 1960 kB      |
pguser | log_ins_201702   | table | pguser | 1768 kB      |
...
pguser | log_ins_201712   | table | pguser | 8192 bytes   |
pguser | log_ins_history   | table | pguser | 1960 kB      |
(14 rows)

```

[illegible]

8.2.4 数据流图

2017-01-01

```
mydb=> EXPLAIN ANALYZE SELECT * FROM log_ins WHERE
create_time > '2017-01-01' AND create_time < '2017-01-02';
QUERY PLAN
```

```
-----
-----
      Append (cost=0.00..45.97 rows=1435 width=16) (actual
time=0.025..0.425 rows=1439 loops=1)
        -> Seq Scan on log_ins (cost=0.00..0.00 rows=1
width=16) (actual time=0.004.. 0.004 rows=0 loops=1)
          Filter: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::times
tamp without time zone))
        -> Index Scan USING idx_log_ins_201701_ctime on
log_ins_201701 (cost=0.29.. 45.97 rows=1434 width=16)
(actual time=0.020..0.285
rows=1439 loops=1)
          Index Cond: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::t
imestamp without time zone))
          Planning time: 0.581 ms
          Execution time: 0.515 ms
(7 rows)
```

```

log_ins_201701
0.515
log_ins_201701
```

```
mydb=> EXPLAIN ANALYZE SELECT * FROM log_ins_201701 WHERE
create_time > '2017-01-01' AND create_time < '2017-01-02';
QUERY PLAN
```

```
-----
-----
      Index Scan USING idx_log_ins_201701_ctime on
log_ins_201701 (cost=0.29..45.97 rows=1434 width=16) (actual
time=0.017..0.254 rows=1
439 loops=1)
        Index Cond: ((create_time > '2017-01-01
```

```

00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::timestamp
        without time zone))
    Planning time: 0.142 ms
    Execution time: 0.337 ms
(4 rows)

```

```

0.337
log_ins

```

8.2.5 constraint_exclusion

```

constraint_exclusion

```

```

·on
·off
·partition UNION ALL

```

```

on partition

```

constraint_exclusion off

```
mydb=> SET constraint_exclusion =off;
SET
```

```
mydb=> EXPLAIN ANALYZE SELECT * FROM log_ins WHERE
create_time > '2017-01-01' AND create_time < '2017-01-02';
QUERY PLAN
-----
Append (cost=0.00..94.40 rows=1447 width=16) (actual
time=0.029..0.534 rows=1439 loops=1)
-> Seq Scan on log_ins (cost=0.00..0.00 rows=1
width=16) (actual time=0.005..0.005 rows=0 loops=1)
Filter: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::timestamp without time zone))
-> Index Scan USING idx_his_ctime on
log_ins_history (cost=0.29..4.31 rows=1 width=16) (actual
time=0.008..0.008 rows=0 loops=1)
Index Cond: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::
timestamp without time zone))
-> Index Scan USING idx_log_ins_201701_ctime on
log_ins_201701 (cost=0.29..45.97 rows=1434 width=16) (actual
time=0.016..0.293
rows=1439 loops=1)
Index Cond: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::
timestamp without time zone))=0 loops=1)
...
-> Seq Scan on log_ins_201712 (cost=0.00..1.01
rows=1 width=16) (actual time=0.011..0.011 rows=0 loops=1)
Filter: ((create_time > '2017-01-01
```

```
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::
        timestamp without time zone))
        Rows Removed by Filter: 1
        Planning time: 1.344 ms
        Execution time: 0.685 ms
(32 rows)
```

0.685
partition
on

8.2.6

log_ins
SQL

```
CREATE TABLE log_ins_201801(CHECK ( create_time >= '2018-01-
01' and create_time < '2018-02-01')) INHERITS(log_ins);
...
```

```
CREATE INDEX idx_log_ins_201801_ctime ON log_ins_201801 USING
btree
(create_time);
...
```

```
CREATE TRIGGER log_ins_insert_trigger
AFTER INSERT ON log_ins_201801
FOR EACH ROW
EXECUTE FUNCTION log_ins_insert_trigger_fn(NEW.log_ins_id, NEW.log_ins_text);
```

```
CREATE TRIGGER log_ins_update_trigger
AFTER UPDATE ON log_ins_201801
FOR EACH ROW
EXECUTE FUNCTION log_ins_update_trigger_fn(NEW.log_ins_id, NEW.log_ins_text);
```

```
-- Create table
CREATE TABLE log_ins_201802(LIKE log_ins INCLUDING ALL );

-- Add constraint
ALTER TABLE log_ins_201802 ADD CONSTRAINT
log_ins_201802_create_time_check
CHECK ( create_time >= '2018-02-01' AND create_time < '2018-
03-01');

-- Create trigger
CREATE TRIGGER log_ins_insert_trigger()
AFTER INSERT ON log_ins_201802
EXECUTE FUNCTION log_ins_insert_trigger_fn(NEW.log_ins_id, NEW.log_ins_text);

-- Inherit table
ALTER TABLE log_ins_201802 INHERIT log_ins;
```

```
CREATE TRIGGER log_ins_delete_trigger
AFTER DELETE ON log_ins_201802
FOR EACH ROW
EXECUTE FUNCTION log_ins_delete_trigger_fn(OLD.log_ins_id, OLD.log_ins_text);
```

8.2.7 删除表

删除表使用 **DELETE** 语句。删除表后，表的所有数据、索引、约束等都将丢失。

```
DROP TABLE log_ins_201802
```

删除表后，表的所有数据、索引、约束等都将丢失。

```
mydb=> ALTER TABLE log_ins_201802 NO INHERIT log_ins;  
ALTER TABLE
```

删除表后，表的所有数据、索引、约束等都将丢失。

8.2.8 删除数据库

删除数据库使用 **\d** 命令。删除数据库后，数据库中的所有数据、索引、约束等都将丢失。

```
mydb=> \d log_ins  
Table "pguser.log_ins"
```

Column	Type	Collation	Nullable	Default
id	integer	not null		
nextval('log_ins_id_seq'::regclass)				
user_id	integer			
create_time	timestamp(0) without time zone			

Triggers:

```
insert_log_ins_trigger BEFORE INSERT ON log_ins FOR EACH
ROW EXECUTE PROCEDURE log_ins_insert_trigger()
```

Number of child tables: 14 (Use \d+ to list them.)

```
\d+log_ins
\d+log_ins_insert_trigger
\d+log_ins
```

```
SQL
```

```
mydb=> SELECT
  nmsp_parent.nspname AS parent_schema ,
  parent.relname AS parent ,
  nmsp_child.nspname AS child_schema ,
  child.relname AS child_schema
FROM
  pg_inherits JOIN pg_class parent
    ON pg_inherits.inhparent = parent.oid JOIN pg_class
child
    ON pg_inherits.inhrelid = child.oid JOIN pg_namespace
nmsp_parent
    ON nmsp_parent.oid = parent.relnamespace JOIN
pg_namespace nmsp_child
    ON nmsp_child.oid = child.relnamespace
WHERE
  parent.relname = 'log_ins';
parent_schema | parent | child_schema | child_schema
-----+-----+-----+-----
pguser        | log_ins | pguser        | log_ins_history
```

pguser	log_ins	pguser	log_ins_201701
pguser	log_ins	pguser	log_ins_201702
...			
pguser	log_ins	pguser	log_ins_201801

(14 rows)

pg_inherits

1. 列出所有表的继承关系
 2. 列出所有表的分区信息
 3. 列出所有表的SQL

```

mydb=> SELECT
    nspname ,
    relname ,
    count(*) AS partition_num
FROM
    pg_class c ,
    pg_namespace n ,
    pg_inherits i
WHERE
    c.oid = i.inhparent
    AND c.relnamespace = n.oid
    AND c.relhassubclass
    AND c.relkind in ('r','p')
GROUP BY 1,2 ORDER BY partition_num DESC;
  nspname | relname | partition_num
-----+-----+-----
  pguser  | log_ins |             14
  pguser  | tbl_log |              1
(2 rows)
  
```

log_ins 14

tbl_log

8.2.9

테이블을 생성하고 로그 테이블에 데이터를 삽입합니다.
create_time은 TIMESTAMP 형식으로, log_ins 테이블의
user_id는 INT 형식으로 지정합니다.

로그 테이블에 데이터를 삽입합니다.
로그 테이블에 데이터를 삽입합니다.

```
CREATE TABLE log(id
serial,user_id int4,
create_time timestamp(0) without time zone);

INSERT INTO log(user_id,create_time)
SELECT round(1000000000*random()),generate_series('2016-12-
01'::date,
'2017-12-01'::date, '1 minute');
```

로그 테이블에 데이터를 삽입합니다.

```
mydb=> SELECT count(*) FROM log_ins;
count
-----
525601
(1 row)
```

```
mydb=> SELECT count(*) FROM log;
count
-----
525601
(1 row)
```

로그 테이블에 데이터를 삽입합니다.

```
CREATE INDEX idx_log_userid ON log USING btree(user_id);
CREATE INDEX idx_log_create_time ON log USING
btree(create_time);
```

log_ins user_id

```
CREATE INDEX idx_log_ins_userid ON log_ins USING
btree(user_id);
CREATE INDEX idx_his_userid ON log_ins_history USING btree
(user_id);
CREATE INDEX idx_log_ins_201701_userid ON log_ins_201701
USING btree (user_id);
CREATE INDEX idx_log_ins_201702_userid ON log_ins_201702
USING btree (user_id);
...
CREATE INDEX idx_log_ins_201801_userid ON log_ins_201801
USING btree (user_id);
```

user_id log_ins

```
-- user_id
SELECT * FROM log WHERE user_id=?;
SELECT * FROM log_ins WHERE user_id=?;
```

log log_ins user_id

```
mydb=> SELECT a.* FROM log a ,log_ins b WHERE
a.user_id=b.user_id LIMIT 1;
   id      | user_id  |      create_time
-----+-----+-----
   67286   | 51751630 | 2017-01-16 17:25:00
(1 row)
```

user_id=51751630log

```
mydb=> EXPLAIN SELECT * FROM log WHERE user_id=51751630;
              QUERY PLAN
-----
Index Scan using idx_log_userid on log  (cost=0.42..4.44
rows=1 width=16)
    Index Cond: (user_id = 51751630)
(2 rows)
```

user_id

```
mydb=> EXPLAIN SELECT * FROM log_ins WHERE user_id=51751630;
              QUERY PLAN
-----
Append  (cost=0.00..63.18 rows=23 width=16)
  -> Seq Scan on log_ins  (cost=0.00..0.00 rows=1
width=16)
    Filter: (user_id = 88258037)
  -> Index Scan USING idx_his_userid on log_ins_history
      (cost=0.29..4.31 rows=1 width=16)
      Index Cond: (user_id = 88258037)
  -> Index Scan USING idx_log_ins_201701_userid on
log_ins_201701
      (cost=0.29.. 4.31 rows=1 width=16)
      Index Cond: (user_id = 88258037)
```

```

-> Index Scan USING idx_log_ins_201702_userid on
log_ins_201702 (cost=0.29.. 4.31 rows=1 width=16)
    Index Cond: (user_id = 88258037)
...
-> Bitmap Heap Scan on log_ins_201801 (cost=2.22..10.48
rows=9 width=16)
    Recheck Cond: (user_id = 88258037)
        -> Bitmap Index Scan on
idx_log_ins_201801_userid (cost=0.00..2.22 rows=9 width=0)
            Index Cond: (user_id = 88258037)
(33 rows)

```

```

log SQL 0.050
log_ins SQL 0.184

create_time log_ins
SQL

```

```

--create_time;
SELECT * FROM log WHERE create_time > '2017-01-01' AND
create_time < '2017-01-02';
SELECT * FROM log_ins WHERE create_time > '2017-01-01' AND
create_time < '2017-01-02';

```

```

log SQL 0.339
log_ins SQL 0.503
8-1

```

8-1

查 询 场 景	普通表 log 执行时间	分区表：查询 log_ins 父表执行时间	分区表：查询 log_ins 子表执行时间
根据非分区键 user_id 查询	0.05 毫秒	0.184 毫秒	不支持
根据分区键 create_time 范围查询	0.339 毫秒	0.503 毫秒	0.325 毫秒

在普通表上，使用非分区键 user_id 查询，执行时间为 2.68 毫秒。在分区表上，使用非分区键 user_id 查询，执行时间为 0.4 毫秒。在分区表上，使用分区键 create_time 查询，执行时间为 0.325 毫秒。

1. 在普通表上，使用非分区键 user_id 查询，执行时间为 2.68 毫秒。

2. 在分区表上，使用非分区键 user_id 查询，执行时间为 0.4 毫秒。

在分区表上，使用分区键 create_time 查询，执行时间为 0.325 毫秒。在普通表上，使用分区键 create_time 查询，执行时间为 0.339 毫秒。在分区表上，使用非分区键 user_id 查询，执行时间为 0.184 毫秒。

8.2.10 索引优化

在普通表上，使用非分区键 user_id 查询，执行时间为 2.68 毫秒。

8.3

```

graph TD
    A[PostgreSQL10] --> B[INSERT]
    B --> C[DELETE]
    C --> D[UPDATE]
    D --> E[DML]
    E --> F[DML]
    F --> G[PostgreSQL10]
  
```

8.3.1 五五五五五

[illegible]

--	--	--	--	--	--	--	--	--

```
CREATE TABLE table_name ( ... )
    [ PARTITION BY { RANGE | LIST } ( { column_name | (
expression ) } ]
```

LIST RANGE

[illegible]

```
CREATE TABLE table_name
    PARTITION OF parent_table [ (
        ) ] FOR VALUES partition bound spec
```

```
partition_bound_spec[]
partition_bound_spec[]
partition_bound_spec[]
partition_bound_spec[]
[]
```

PostgreSQL10

```
1
2
3
log_par
```

```
CREATE TABLE log_par (
  id serial,
  user_id int4,
  create_time timestamp(0) without time zone
) PARTITION BY RANGE(create_time);
```

```
log_par
create_time
```


□□□□□□□□□□□□□□

```
INSERT INTO log_par(user_id,create_time)
SELECT round(100000000*random()),generate_series('2016-12-
01'::date,
          '2017-12-01'::date, '1 minute');
```

□□□□□□□□□□□□□□

```
mydb=> SELECT count(*) FROM log_par;
count
-----
525601
(1 row)

mydb=> SELECT count(*) FROM ONLY log_par;
count
-----
0
(1 row)
```

□□□□□□□□□□□□log_par□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
mydb=> \dt+ log_par*

          List of relations
 Schema |      Name      | Type  | Owner  |      Size      |
-----+-----+-----+-----+-----+-----+
 pguser | log_par        | table | pguser | 0 bytes        |
 pguser | log_par_201701 | table | pguser | 1960 kB        |
 pguser | log_par_201702 | table | pguser | 1768 kB        |
 ...
```

pguser		log_par_201712		table		pguser		8192 bytes	
pguser		log_par_his		table		pguser		1960 kB	

8.3.3 数据库链接

数据库链接是指数据库与数据库之间的连接。在 PostgreSQL 中，可以通过 `log_par` 表来记录数据库链接的信息。

```
mydb=> SELECT
    nmsp_parent.nspname AS parent_schema ,
    parent.relname AS parent ,
    nmsp_child.nspname AS child_schema ,
    child.relname AS child_schema
FROM
    pg_inherits JOIN pg_class parent
        ON pg_inherits.inhparent = parent.oid JOIN pg_class
child
        ON pg_inherits.inhrelid = child.oid JOIN
pg_namespace nmsp_parent
        ON nmsp_parent.oid = parent.relnamespace JOIN
pg_namespace nmsp_child
        ON nmsp_child.oid = child.relnamespace
WHERE
    parent.relname = 'log_par';
parent_schema | parent | child_schema | child_schema
-----+-----+-----+-----
pguser        | log_par | pguser        | log_par_his
pguser        | log_par | pguser        | log_par_201701
pguser        | log_par | pguser        | log_par_201702
...
pguser        | log_par | pguser        | log_par_201712
(13 rows)
```

SQL 语句用于查询 `log_par` 表中的记录。例如，可以使用以下语句来查询所有以 `log_par` 开头的记录：

```
\d+log_par
log_par
log_par
```

8.3.4 分区

创建分区表 log_par

```
CREATE TABLE log_par_201801 PARTITION OF log_par FOR VALUES  
FROM ('2018-01-01') TO ('2018-02-01');
```

创建索引

```
CREATE INDEX idx_log_par_201801_ctime ON log_par_201801  
USING btree(create_time);
```

8.3.5 删除

删除分区 log_par_201801

```
DROP TABLE log_par_201801;
```

删除表 log_par

```
mydb=> ALTER TABLE log_par DETACH PARTITION log_par_201801;  
ALTER TABLE
```


log_par_201701
0.510

log_par log
create_time log_par user_id
SQL

log_par user_id

```
CREATE INDEX idx_log_par_his_userid ON log_par_his using  
btree (user_id);  
CREATE INDEX idx_log_par_201701_userid ON log_par_201701  
using btree (user_id);  
CREATE INDEX idx_log_par_201702_userid ON log_par_201702  
using btree (user_id);  
...  
CREATE INDEX idx_log_par_201712_userid ON log_par_201712  
using btree (user_id);
```

user_id log_par
SQL

```
-- user_id  
SELECT * FROM log WHERE user_id=?;  
SELECT * FROM log_par WHERE user_id=?;
```

log log_par user_id

```
mydb=> SELECT a.* FROM log a ,log_par b WHERE
a.user_id=b.user_id LIMIT 1;
      id | user_id |      create_time
-----+-----+-----
    258926 | 70971018 | 2017-05-29 19:25:00
(1 row)
```

user_id=70971018log

```
mydb=> EXPLAIN SELECT * FROM log WHERE user_id=70971018;
      QUERY PLAN
-----
Index Scan using idx_log_userid on log (cost=0.42..4.44
rows=1 width=16)
    Index Cond: (user_id = 70971018)
(2 rows)
```

user_idlog_par

```
mydb=> EXPLAIN SELECT * FROM log_par WHERE user_id=70971018;
      QUERY PLAN
-----
Append (cost=0.29..52.70 rows=13 width=16)
-> Index Scan using idx_log_par_his_userid on
log_par_his (cost=0.29..4.31 rows=1 width=16)
    Index Cond: (user_id = 70971018)
-> Index Scan using idx_log_par_201701_userid on
log_par_201701
```

```

      (cost=0.29.. 4.31 rows=1 width=16)
      Index Cond: (user_id = 70971018)
-> Index Scan using idx_log_par_201702_userid on
log_par_201702
      (cost=0.29.. 4.31 rows=1 width=16)
      Index Cond: (user_id = 70971018)
...
-> Seq Scan on log_par_201712 (cost=0.00..1.01 rows=1
width=16)
      Filter: (user_id = 70971018)
(27 rows)

```

```

      user_id
log_par      log      SQL
0.047log_par      SQL
0.139

```

```

      create_time      log_par
SQL

```

```

--create_time;
SELECT * FROM log WHERE create_time > '2017-01-01' AND
create_time < '2017-01-02';
SELECT * FROM log_par WHERE create_time > '2017-01-01' AND
create_time < '2017-01-02';

```

```

log      SQL0.340
log_par      SQL0.503
8-2

```

8-2

查 询 场 景	普通表 log 执行时间	分区表：查询 log_par 父表执行时间	区表：查询 log_par 子表执行时间
根据非分区键 user_id 查询	0.047 毫秒	0.139 毫秒	不支持
根据分区键 create_time 范围查询	0.340 毫秒	0.503 毫秒	0.319 毫秒

1.95
 create_time
 0.47
 create_time
 user_id
 8-1
 8-2
 8-3

8-3

查 询 场 景	普通表 (log)	传统分区表 (log_ins)		内置分区表 (log_par)	
		查询父表	查询子表	查询父表	查询子表
根据非分区键 user_id 查询	表 8-1：0.05 毫秒 表 8-2：0.047 毫秒	0.184 毫秒	不支持	0.139 毫秒	不支持
根据分区键 create_time 范围查询	表 8-1：0.339 毫秒 表 8-2：0.340 毫秒	0.503 毫秒	0.325 毫秒	0.503 毫秒	0.319 毫秒

•

•


```

-> Index Scan using idx_log_par_201701_ctime on
log_par_201701 (cost=0.29..45.21 rows=1396 width=16)
      (actual time=0.016..0.280 rows=1439 loops=1)
    Index Cond: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND (create_time <
'2017-01-02 00:00:00'::timestamp without time zone))
...
-> Bitmap Index Scan on
idx_log_par_201801_ctime (cost=0.00..2.24 rows=9 width=0)
      (actual time=0.002..0.002 rows=0 loops=1)
    Index Cond: ((create_time > '2017-01-01
00:00:00'::timestamp without time zone) AND
      (create_time < '2017-01-02
00:00:00'::timestamp without time zone))
    Planning time: 0.792 ms
    Execution time: 0.607 ms
(34 rows)

```

0.607 partition on

8.3.8

UPDATE SQL

```

mydb=> SELECT * FROM log_par_201701 LIMIT 1;
   id  | user_id |   create_time
-----+-----
 44641 | 16965492 | 2017-01-01 00:00:00
(1 row)

```

```

mydb=> UPDATE log_par SET create_time='2017-02-02 01:01:01'

```

```
WHERE user_id=16965492;
ERROR:  new row for relation "log_par_201701" violates
partition constraint
DETAIL:  Failing row contains (44641, 16965492, 2017-02-02
01:01:01).
```

```
user_id16965492
log_par_201701create_time
'2017-02-0201:01:01'
```

```
mydb=> UPDATE log_par SET create_time='2017-01-01 01:01:01'
WHERE user_id=16965492;
UPDATE 1
```

8.3.9

```
•
```

```
•
```


8.4 数据库

数据库是指长期存储在计算机内、有组织的、可共享的数据集合。数据库中的数据按一定的数据模型组织、存储于数据库中。数据库系统是指引入数据库技术后，对数据库进行管理的软件系统。数据库系统由数据库、数据库管理系统（DBMS）、数据库用户等组成。数据库系统的主要功能包括：数据定义、数据操纵、数据控制、数据查询、数据维护等。数据库系统的发展经历了多个阶段，从早期的文件系统到现在的关系数据库、NoSQL数据库等。数据库系统是现代信息系统的核心组成部分，广泛应用于各个领域。

PostgreSQL10 是一款开源的关系数据库管理系统，支持 SQL 标准，具有强大的功能和良好的性能。它支持多种数据类型，包括文本、数字、日期、时间、几何等。PostgreSQL10 还支持事务处理、并发控制、备份恢复等功能。它广泛应用于企业级应用、大数据分析、科学计算等领域。PostgreSQL10 的官方网站为 <https://www.postgresql.org/>。

SQL 是一种用于管理和操作数据库的查询语言。它包括数据定义语言（DDL）、数据操纵语言（DML）、数据控制语言（DCL）等。SQL 语言简单易学，功能强大，是数据库开发人员必须掌握的技能。SQL 语言的发展经历了多个阶段，从早期的结构化查询语言（SQL）到现在的 SQL:2011 标准。SQL 语言广泛应用于各种数据库系统，包括关系数据库、NoSQL 数据库等。

9 PostgreSQL vs NoSQL

PostgreSQL

3

PostgreSQL json jsonb

json jsonb jsonb json

jsonb json jsonb

9.1 jsonb数据类型

jsonb数据类型是json类型的二进制表示，它比json类型更高效，并且支持索引。GIN索引是用于jsonb数据类型的一种索引类型。

```
{
  "id": 1,
  "user_id": 1440933,
  "user_name": "l_francs",
  "create_time": "2017-08-03 16:22:05.528432+08"
}
```

jsonb数据类型可以用于存储JSON数据。例如，我们可以创建一个表，其中包含一个jsonb类型的字段，用于存储用户信息。

```
CREATE INDEX idx_gin ON tbl_user_jsonb USING gin(user_info);
```

jsonb数据类型支持GIN索引。我们可以使用@>操作符来查询jsonb数据中的特定字段。

```
SELECT * FROM tbl_user_jsonb WHERE user_info @>
'{"user_name": "l_francs"}'
```

jsonb数据类型支持GIN索引。我们可以使用idx_gin索引来加速查询。

```
SELECT * FROM tbl_user_jsonb WHERE user_info->>'user_name'='1_francs';
```

tbl_user_jsonb tbl_user_jsonb
tbl_user_jsonb

```
CREATE INDEX idx_gin_user_info_b_user_name ON tbl_user_jsonb  
USING btree  
((user_info ->> 'user_name'));
```

tbl_user_info->>'user_name'
>>'user_name' SQL

9.2 json/jsonb

jsonb는 json과 유사하지만, json과 달리
jsonb는 json과 달리 3가지 특징이 있다.
1. jsonb는 json과 달리, json과 달리
2. jsonb는 json과 달리, json과 달리
3. jsonb는 json과 달리, json과 달리
4. jsonb는 json과 달리, json과 달리
5. jsonb는 json과 달리, json과 달리
6. jsonb는 json과 달리, json과 달리
7. jsonb는 json과 달리, json과 달리
8. jsonb는 json과 달리, json과 달리
9. jsonb는 json과 달리, json과 달리
10. jsonb는 json과 달리, json과 달리

9.2.1 json/jsonb

json/jsonb는 json/jsonb와 유사하지만,
json/jsonb는 json/jsonb와 유사하지만,
json/jsonb는 json/jsonb와 유사하지만,

·user_ini int(4) unsigned zerofill(200)

·tbl_user_json json(200)

·tbl_user_jsonb jsonb(200)

·user_ini int(4) unsigned zerofill(200)

```
mydb=> CREATE TABLE user_ini(id int4 ,user_id int8, user_name  
character
```

```
varying(64),create_time timestamp(6) with time zone default  
clock_timestamp());  
CREATE TABLE
```

```
mydb=> INSERT INTO user_ini(id,user_id,user_name)  
SELECT r,round(random()*2000000), r || '_francs'  
FROM generate_series(1,2000000) as r;  
INSERT 0 2000000
```

user_ini json jsonb
user_ini_json user_ini_jsonb

```
mydb=> CREATE TABLE tbl_user_json(id serial, user_info json);  
CREATE TABLE  
mydb=> CREATE TABLE tbl_user_jsonb(id serial, user_info  
jsonb);  
CREATE TABLE
```

9.2.2 json jsonb

user_ini row_to_json
user_ini_json 200 json

```
mydb=> \timing  
Timing is on.
```

```
mydb=> INSERT INTO tbl_user_json(user_info) SELECT  
row_to_json(user_ini)  
FROM user_ini;  
INSERT 0 2000000  
Time: 13825.974 ms (00:13.826)
```

tbl_user_json 200 13
user_ini 200 jsonb
tbl_user_jsonb

```
mydb=> INSERT INTO tbl_user_jsonb(user_info)
        SELECT row_to_json(user_ini)::jsonb FROM user_ini;
INSERT 0 2000000
Time: 20756.993 ms (00:20.757)
```

tbl_user_jsonb 200 jsonb
20 json jsonb
tbl_user_jsonb

```
mydb=> \dt+ tbl_user_json
          List of relations
 Schema |      Name      | Type  | Owner  | Size  |
Description
-----+-----+-----+-----+-----+-----
--
 pguser | tbl_user_json | table | pguser | 281 MB |
(1 row)

mydb=> \dt+ tbl_user_jsonb
          List of relations
 Schema |      Name      | Type  | Owner  | Size  |
Description
-----+-----+-----+-----+-----+-----
---
 pguser | tbl_user_jsonb | table | pguser | 333 MB |
(1 row)
```

jsonb json
tbl

tbl-user-json

```
mydb=> SELECT * FROM tbl_user_json LIMIT 1;
      id      |
 user_info    |
-----+-----
      2000001 |
 {"id":1,"user_id":1182883,"user_name":"l_francs","create_time"
 : "2017-08-03T20:59:27.42741+08:00"}
 (1 row)
```

9.2.3 json/jsonb

json/jsonb 数据类型支持 json/jsonb 格式的数据存储。在 user_info 表中，user_name 字段为 json/jsonb 类型。

```
mydb=> EXPLAIN ANALYZE SELECT * FROM tbl_user_jsonb WHERE
 user_info->>'user_name'='l_francs';
QUERY PLAN
-----
Seq Scan on tbl_user_jsonb  (cost=0.00..72859.90 rows=10042
width=143) (actual time=0.023..524.843 rows=1 loops=1)
  Filter: ((user_info ->> 'user_name'::text) =
'l_francs'::text)
  Rows Removed by Filter: 1999999
  Planning time: 0.091 ms
  Execution time: 524.876 ms
(5 rows)
```

SQL 语句中 524 毫秒的时间主要消耗在 user_info 表的 user_name 字段的 btree 索引扫描上。

```
mydb=> CREATE INDEX idx_jsonb ON tbl_user_jsonb USING btree
((user_info->>'user_name'));
```

□□□□□□□□□□□□□□□□

```
mydb=> EXPLAIN ANALYZE SELECT * FROM tbl_user_jsonb WHERE
user_info->>'user_name'='1_francs';
```

QUERY PLAN

```
-----
-----
      Bitmap Heap Scan on tbl_user_jsonb (cost=155.93..14113.93
rows=10000 width=143) (actual time=0.027..0.027 rows=1
loops=1)
    Recheck Cond: ((user_info ->> 'user_name'::text) =
'1_francs'::text)
    Heap Blocks: exact=1
    -> Bitmap Index Scan on idx_jsonb (cost=0.00..153.43
rows=10000 width=0) (actual time=0.021..0.021 rows=1 loops=1)
        Index Cond: ((user_info ->> 'user_name'::text) =
'1_francs'::text)
    Planning time: 0.091 ms
    Execution time: 0.060 ms
(7 rows)
```

□□□□□□□□□□□□□□□□□□□□SQL□□□□□□
0.060ms□□□□□□□□tbl_user_jsonb□
tbl_user_jsonb□□□□□□□□□□□□□□□□user_info□□
id□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
mydb=> CREATE INDEX idx_gin_user_info_id ON tbl_user_json
USING btree
(((user_info ->> 'id')::integer));
CREATE INDEX
```

```
mydb=> CREATE INDEX idx_gin_user_infob_id ON tbl_user_jsonb
USING btree
```

```
((user_info -> 'id')::integer));  
CREATE INDEX
```

tbl_user_json

```
mydb=> EXPLAIN ANALYZE SELECT id,user_info->'id',user_info->'user_name'  
FROM tbl_user_json  
WHERE (user_info->'id')::int4>1 AND (user_info->'id')::int4<10000;
```

QUERY PLAN

```
-----  
-----  
Bitmap Heap Scan on tbl_user_json (cost=166.30..14178.17  
rows=10329 width=68) (actual time=1.167..26.534 rows=9998  
loops=1)  
Recheck Cond: (((user_info -> 'id')::text))::integer  
> 1) AND (((user_info -> 'id')::text))::integer < 10000)  
Heap Blocks: exact=338  
-> Bitmap Index Scan on idx_gin_user_info_id  
(cost=0.00..163.72 rows=10329 width=0) (actual  
time=1.110..1.110 rows=19996 loops=1)  
Index Cond: (((user_info ->  
'id')::text))::integer > 1) AND (((user_info ->  
'id')::text))::integer < 10000))  
Planning time: 0.094 ms  
Execution time: 27.092 ms  
(7 rows)
```

tbl_user_json user_info
id 1 10000
27.092 tbl_user_jsonb SQL

```
mydb=> EXPLAIN ANALYZE SELECT id,user_info->'id',user_info->'user_name'  
FROM tbl_user_jsonb
```

```
WHERE (user_info->>'id')::int4>1 AND (user_info->>'id')::int4<10000;
```

```
QUERY PLAN
```

```
-----  
-----  
      Bitmap Heap Scan on tbl_user_jsonb (cost=158.93..14316.93  
rows=10000 width=68) (actual time=1.140..8.116 rows=9998  
loops=1)  
        Recheck Cond: (((user_info ->> 'id')::text))::integer  
> 1) AND (((user_info ->> 'id')::text))::integer < 10000)  
        Heap Blocks: exact=393  
        -> Bitmap Index Scan on idx_gin_user_infob_id  
(cost=0.00..156.43 rows= 10000 width=0) (actual  
time=1.058..1.058 rows=18992 loops=1)  
          Index Cond: (((user_info ->>  
'id')::text))::integer > 1) AND (((user_info ->>  
'id')::text))::integer < 10000)  
          Planning time: 0.104 ms  
          Execution time: 8.656 ms  
(7 rows)
```

tbl_user_jsonb
user_info[id]110000
8.656jsonbjson

“jsonjsonb
jsonb”key/value

```
SELECT * FROM tbl_user_jsonb WHERE user_info @> '{"user_name":  
"2_francs"}';
```

```
mydb=> EXPLAIN ANALYZE SELECT * FROM tbl_user_jsonb WHERE
user_info @> '{"user_name": "2_francs"}';
          QUERY PLAN
```

```
-----
Seq Scan on tbl_user_jsonb (cost=0.00..67733.00 rows=2000
width=143) (actual time=0.018..582.207 rows=1 loops=1)
  Filter: (user_info @> '{"user_name":
"2_francs"}'::jsonb)
  Rows Removed by Filter: 1999999
  Planning time: 0.065 ms
  Execution time: 582.232 ms
(5 rows)
```

582ms
tbl_user_jsonb user_info gin

```
mydb=> CREATE INDEX idx_tbl_user_jsonb_user_Info ON
tbl_user_jsonb USING gin
      (user_Info);
CREATE INDEX
```

```
mydb=> EXPLAIN ANALYZE SELECT * FROM tbl_user_jsonb WHERE
user_info @> '{"user_name": "2_francs"}';
          QUERY PLAN
```

```
-----
Bitmap Heap Scan on tbl_user_jsonb (cost=37.50..3554.34
rows=2000 width=143) (actual time=0.079..0.080 rows=1 loops=1)
  Recheck Cond: (user_info @> '{"user_name":
"2_francs"}'::jsonb)
  Heap Blocks: exact=1
  -> Bitmap Index Scan on idx_tbl_user_jsonb_user_info
(cost=0.00..37.00 rows=2000 width=0) (actual time=0.069..0.069
rows=1 loops=1)
    Index Cond: (user_info @> '{"user_name":
```

```
"2_francs"}'::jsonb)
  Planning time: 0.094 ms
  Execution time: 0.114 ms
(7 rows)
```

0.114

json jsonb json
jsonb jsonb

9.3 json/jsonb数据类型

json/jsonb数据类型是PostgreSQL10引入的，用于存储JSON数据。json数据类型是文本类型，而jsonb数据类型是二进制类型，支持索引和查询。PostgreSQL支持json和jsonb两种数据类型，其中jsonb是更推荐使用的。

9.3.1 PostgreSQL中的JSON

PostgreSQL中的JSON数据类型支持两种格式：JSON和JSONB。JSON是文本格式，而JSONB是二进制格式。JSONB格式支持索引和查询，而JSON格式不支持。PostgreSQL中的JSON数据类型可以用于存储和查询JSON数据。

在PostgreSQL中，可以使用`like`操作符来查询JSON数据。例如，可以使用`like`操作符来查询JSON数据中的特定字段。

- 使用`like`操作符来查询JSON数据中的特定字段。

- 使用`like`操作符来查询JSON数据中的特定字段。

PostgreSQL中的JSON数据类型可以用于存储和查询JSON数据。PostgreSQL中的JSON数据类型可以用于存储和查询JSON数据。

1.tsvector

tsvector는 텍스트를 토큰화하고, 각 토큰에 대해 역색인 정보를 제공하는 데이터 타입입니다. tsvector는 텍스트를 토큰화하고, 각 토큰에 대해 역색인 정보를 제공하는 데이터 타입입니다.

```
mydb=> SELECT 'Hello,cat,how are u? cat is smiling!'
::tsvector;
          tsvector
-----
'Hello,cat,how' 'are' 'cat' 'is' 'smiling!' 'u?'
(1 row)
```

to_tsvector는 텍스트를 토큰화하고, 각 토큰에 대해 역색인 정보를 제공하는 데이터 타입입니다. to_tsvector는 텍스트를 토큰화하고, 각 토큰에 대해 역색인 정보를 제공하는 데이터 타입입니다.

```
mydb=> SELECT to_tsvector('english','Hello cat,');
          to_tsvector
-----
'cat':2 'hello':1
(1 row)
```

2.tsquery

tsquery는 텍스트 쿼리를 저장하는 데이터 타입입니다. tsquery는 텍스트 쿼리를 저장하는 데이터 타입입니다.

```
mydb=> SELECT 'hello&cat'::tsquery;
          tsquery
-----
```



```
'hello' & 'cat'
(1 row)
```

to_tsquery

```
mydb=> SELECT to_tsquery( 'hello&cat' );
          to_tsquery
-----
'hello' & 'cat'
(1 row)
```

to_tsvector

```
mydb=> SELECT to_tsvector('english','Hello cat,how are u')
@@
to_tsquery( 'hello&cat' );
?column?
-----
t
(1 row)
```

to_tsvector

```
mydb=> SELECT to_tsvector('english','Hello cat,how are u')
@@
to_tsquery( 'hello&dog' );
?column?
-----
```

f
(1 row)

tsquery"|"



to_tsvector
to_tsvector
[config regconfig]document text
to_tsvectorconfigenglish
configdefault_text_search_config

3.

200

```
mydb=> CREATE TABLE test_search(id int4,name text);
CREATE TABLE
mydb=> INSERT INTO test_search(id,name) SELECT n,
n||'_francs'
FROM generate_series(1,2000000) n;
INSERT 0 2000000
```

SQLtest_searchname
1_francs

```
mydb=> SELECT * FROM test_search WHERE name LIKE '1_francs';
   id |   name
-----+-----
    1 | 1_francs
(1 row)
```

□□□□□□□□

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_search WHERE name
LIKE '1_francs';
```

QUERY PLAN

```
-----
Seq Scan on test_search (cost=0.00..38465.04 rows=204
width=18) (actual time=0.022..261.766 rows=1 loops=1)
  Filter: (name ~~ '1_francs'::text)
  Rows Removed by Filter: 1999999
  Planning time: 0.101 ms
  Execution time: 261.796 ms
(5 rows)
```

□□□□□□□□□□□□□□□□□□□□261□□□□□□□□□□
□□□□□□□□□□□□□□

```
mydb=> CREATE INDEX idx_gin_search ON test_search USING gin
(to_tsvector('english',name));
CREATE INDEX
```

□□□□SQL□□□test_search□name□□□□□□
1_francs□□□□

```
mydb=> SELECT * FROM test_search WHERE
to_tsvector('english',name) @@
to_tsquery('english','1_francs');
   id |   name
-----+-----
    1 | 1_francs
(1 row)
```

□□□□□□□□□□□□□□□□□□□□

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_search WHERE
to_tsvector('english', name) @@
to_tsquery('english','1_francs');
QUERY PLAN
-----
      Bitmap Heap Scan on test_search  (cost=18.39..128.38
rows=50 width=36) (actual time=0.071..0.071 rows=1 loops=1)
        Recheck Cond: (to_tsvector('english'::regconfig,
name) @@ '''1' & 'franc'::tsquery)
        Heap Blocks: exact=1
        -> Bitmap Index Scan on idx_gin_search
        (cost=0.00..18.38 rows=50 width=0) (actual time=0.064..0.064
rows=1 loops=1)
          Index Cond: (to_tsvector('english'::regconfig,
name) @@ '''1' & 'franc'::tsquery)
        Planning time: 0.122 ms
        Execution time: 0.104 ms
(7 rows)
```

□□□□□□□□□□□□□□□□□□□□0.104□□□□
□□□□3□□□□□□□□□□□□□□SQL□□□□□□□□□□□□

```
mydb=> EXPLAIN ANALYZE SELECT * FROM test_search
WHERE to_tsvector(name) @@ to_tsquery('1_francs');
QUERY PLAN
-----
```

```

-----
Seq Scan on test_search (cost=0.00..1037730.00 rows=50
width=18) (actual time=0.036..10297.764 rows=1 loops=1)
  Filter: (to_tsvector(name) @@
to_tsquery('1_francs'::text))
  Rows Removed by Filter: 1999999
  Planning time: 0.098 ms
  Execution time: 10297.787 ms
(5 rows)

```

```

to_tsvector('english'
name where
to_tsvector to_tsvector
name

```

9.3.2 json jsonb

PostgreSQL 10 json jsonb
 10 json jsonb
 10

1. PostgreSQL 10 9.6 to_tsvector

9.6 to_tsvector

```

[postgres@pghost1 ~]$ psql francs francs
psql (9.6.3)
Type "help" for help.

```

```

mydb=> \df *to_tsvector*

```

Schema	Name	List of functions Result data type
--------	------	---

Argument data types		Type	
-----+-----		-----+-----	
pg_catalog	array_to_tsvector	tsvector	text[]
normal			
pg_catalog	to_tsvector	tsvector	
regconfig, text	normal		
pg_catalog	to_tsvector	tsvector	text
normal			
(3 rows)			

9.6 to_tsvector
 text[]
 10 to_tsvector

```

[postgres@pghost1 ~]$ psql mydb pguser
psql (10.0)
Type "help" for help.
mydb=> \df *to_tsvector*

```

Schema		Name	List of functions	
Argument data types		Type	Result data type	
-----+-----		-----+-----	-----+-----	
pg_catalog	array_to_tsvector	tsvector		text[]
normal				
pg_catalog	to_tsvector	tsvector		json
normal				
pg_catalog	to_tsvector	tsvector		jsonb
normal				
pg_catalog	to_tsvector	tsvector		
regconfig, json	normal			
pg_catalog	to_tsvector	tsvector		
regconfig, jsonb	normal			
pg_catalog	to_tsvector	tsvector		
regconfig, text	normal			
pg_catalog	to_tsvector	tsvector		text
normal				
(7 rows)				

10to_tsvector
jsonjsonb

2.

random_rangeint4int4

```
CREATE OR REPLACE FUNCTION random_range(int4, int4)
RETURNS int4
LANGUAGE SQL
AS $$
    SELECT ($1 + FLOOR(($2 - $1 + 1) * random()))::int4;
$$;
```

random_text_simplelength int4
random_rangeint4int4

```
CREATE OR REPLACE FUNCTION random_text_simple(length int4)
RETURNS text
LANGUAGE PLPGSQL
AS $$
DECLARE
    possible_chars text :=
'0123456789ABCDEFGHIJKLMNOPQRSTUVWXYZ';
    output text := '';
    i int4;
    pos int4;
BEGIN
    FOR i IN 1..length LOOP
        pos := random_range(1, length(possible_chars));
```

```

        output := output || substr(possible_chars, pos, 1);
    END LOOP;

    RETURN output;
END;
$$;

```

random_text_simple(length int4)

```

mydb=> SELECT random_text_simple(3);
       random_text_simple
-----
      LL9
(1 row)

```

```

mydb=> SELECT random_text_simple(6);
       random_text_simple
-----
     B81BPW
(1 row)

```

3. json

user_ini random_text_simple
length int4 100

□□□□□□□□□□

```
mydb=> CREATE TABLE user_ini(id int4 ,user_id int8,
user_name character varying(64),
create_time timestamp(6) with time zone default
clock_timestamp());
CREATE TABLE

mydb=> INSERT INTO user_ini(id,user_id,user_name)
SELECT r,round(random()*1000000), random_text_simple(6)
FROM generate_series(1,1000000) as r;
INSERT 0 1000000
```

□□tbl_user_search_json□□□□□
row_to_json□□□□user_ini□□□□□□□json□□□□□
□□□

```
mydb=> CREATE TABLE tbl_user_search_json(id serial,
user_info json);
CREATE TABLE

mydb=> INSERT INTO tbl_user_search_json(user_info)
SELECT row_to_json(user_ini) FROM user_ini;
INSERT 0 1000000
```

□□□□□□□□□□

```
mydb=> SELECT * FROM tbl_user_search_json LIMIT 1;
      id      |
user_info
-----+-----
          1 |
{"id":1,"user_id":186536,"user_name":"KTU89H","create_time":
```

```
"2017-08-05T15:59:25.359148+08:00"}
(1 row)
```

4.json테이블

테이블 tbl_user_search_json
user_info에 KTU89H에 대한 검색 결과

```
mydb=> SELECT * FROM tbl_user_search_json
WHERE to_tsvector('english',user_info) @@
to_tsquery('ENGLISH','KTU89H');
   id   |
user_info
-----+-----
      1 |
{"id":1,"user_id":186536,"user_name":"KTU89H","create_time":
"2017-08-05T15:59:25.359148+08:00"}
(1 row)
```

SQL테이블을 json테이블로
SQL테이블로 변환하는 방법 8061

```
mydb=> EXPLAIN ANALYZE SELECT * FROM tbl_user_search_json
      WHERE to_tsvector('english',user_info) @@
to_tsquery('ENGLISH','KTU89H');
QUERY PLAN
-----
Seq Scan on tbl_user_search_json  (cost=0.00..279513.00
rows=5000 width=104) (actual time=0.046..8061.858 rows=1
loops=1)
   Filter: (to_tsvector('english'::regconfig,
user_info) @@ ''ktu89h''::tsquery)
   Rows Removed by Filter: 999999
```

```
Planning time: 0.091 ms
Execution time: 8061.880 ms
(5 rows)
```

□□□□□□□

```
mydb=> CREATE INDEX idx_gin_search_json ON
tbl_user_search_json USING
gin(to_tsvector('english',user_info));
CREATE INDEX
```

□□□□□□□□□□SQL□□□□□□

```
mydb=> EXPLAIN ANALYZE SELECT * FROM tbl_user_search_json
WHERE to_tsvector('english', user_info) @@
to_tsquery('ENGLISH','KTU89H');
QUERY PLAN
```


```
      Bitmap Heap Scan on tbl_user_search_json
(cost=50.75..7876.06 rows=5000 width=104) (actual
time=0.024..0.024 rows=1 loops=1)
    Recheck Cond: (to_tsvector('english'::regconfig,
user_info) @@ '''ktu89h''':: tsquery)
    Heap Blocks: exact=1
    -> Bitmap Index Scan on idx_gin_search_json
(cost=0.00..49.50 rows=5000 width=0) (actual
time=0.018..0.018 rows=1 loops=1)
        Index Cond: (to_tsvector('english'::regconfig,
user_info) @@ '''ktu89h'''::
        tsquery)
Planning time: 0.113 ms
Execution time: 0.057 ms
(7 rows)
```

0.057

PostgreSQL
PostgreSQL10
json jsonb

9.4 数据库

数据库数据库PostgreSQLNoSQL数据库数据库
jsonb数据库数据库数据库数据库数据库jsonjsonb数据库
json数据库数据库数据库数据库数据库PostgreSQL数据库数据库
数据库数据库jsonjsonb数据库数据库PostgreSQL10数据库
数据库数据库数据库数据库jsonjsonb数据库数据库数据库数据库数据库
数据库PostgreSQL数据库数据库数据库数据库数据库数据库数据库
PostgreSQL数据库数据库数据库数据库数据库数据库数据库数据库数据库数据库

□□□

- 10□ □□□□
- 11□ □□□□pgbench
- 12□ □□□□□□□□
- 13□ □□□□□
- 14□ □□□
- 15□ □□□□
- 16□ □□□□
- 17□ Oracle□□□□□PostgreSQL□□
- 18□ PostGIS

10

The diagram illustrates a data center layout with a grid of 100 small squares. The layout is divided into four main sections:

- CPU**: Located in the top right corner, consisting of 10 squares.
- I/O**: Located in the top left corner, consisting of 10 squares.
- PostgreSQL**: Located in the bottom left corner, consisting of 10 squares.
- SQL**: Located in the bottom right corner, consisting of 10 squares.

The remaining 70 squares in the grid are empty, representing other potential server locations or unused space.

PostgreSQL

PostgreSQL

SQL

OLTP

10.1 五五五五五

□□□□□□□□□□□□□□CPU□□□□□□□□□□□□

```

graph TD
    IO[I/O] --> SATA[SATA SSD]
    IO --> PCIe[PCIe SSD]
    IO --> SSD[SSD]
    SSD --> SAN[SAN]
    SSD --> NAS[NAS]
    NAS --> NFS[NFS]
  
```

The diagram illustrates a sequence of components in a system boot process. It consists of a horizontal row of rectangles. From left to right, the components are:

- A large rectangle labeled "CPU".
- A smaller rectangle labeled "CPU".
- A rectangle labeled "PostgreSQL".
- A rectangle labeled "BIOS".
- A rectangle labeled "CPU".
- A rectangle labeled "CPU".

The rectangles are arranged in a way that suggests a flow or sequence, with some components being larger than others, possibly indicating their relative importance or duration in the process.

The diagram illustrates a computer system architecture. At the center is a large rectangle labeled "CPU". To the right of the CPU is a large rectangle labeled "I/O". Below the CPU, there are several smaller rectangles representing different types of I/O devices: a keyboard, a mouse, a monitor, a printer, and a scanner. Arrows indicate the flow of data between the CPU and these devices.

10.2 性能调优

性能调优是指通过调整系统参数、配置、硬件等方式，提高系统的运行效率和资源利用率。DBA 需要了解系统的性能瓶颈，并采取相应的措施进行优化。性能调优的目标是在保证系统稳定性和可靠性的前提下，提高系统的吞吐量和响应时间。

10.2.1 Linux 性能调优

Linux 性能调优涉及多个方面，包括 CPU、内存、磁盘 I/O、网络等。常用的性能调优工具包括 top、free、vmstat、iostat、mpstat、sar、pidstat 等。这些工具可以帮助 DBA 了解系统的运行状态，并找出性能瓶颈。

CentOS 安装 sysstat 的步骤如下：

```
[root@pghost1 ~]# yum install -y sysstat
```

1.top

top 是一个常用的性能调优工具，可以实时显示系统的运行状态。通过 top 命令，DBA 可以查看系统的 CPU 使用率、内存使用情况、进程列表等信息。top 命令的输出结果如下：

top

```
top - 12:01:03 up 93 days, 23:30,  1 user,  load average:
1.08, 1.09, 1.08
Tasks: 1042 total,   3 running, 1039 sleeping,   0 stopped,
0 zombie
Cpu(s):  5.3%us,   1.6%sy,   0.0%ni, 92.5%id,   0.2%wa,
0.0%hi,   0.3%si,   0.0%st
Mem:  330604220k total, 323235920k used,  7368300k free,
248996k buffers
Swap: 67108860k total,      0k used, 67108860k free,
295053464k cached

      PID USER      PR  NI  VIRT  RES  SHR  S  %CPU  %MEM
TIME+  COMMAND
145138 postgres  20   0 37.7g  31g  31g  S  26.2  10.1
5:38.22 postgres: pguser mydb 127.0.0.1(50382)
      58327 pgbounce  20   0 49368 8808  864  S  17.4   0.0
14725:14 /usr/bin/pgbouncer -d -q
/etc/pgbouncer/pgbouncer.ini
183682 postgres  20   0 37.6g  22g  21g  R  17.3   7.0
1:49.69 postgres: pguser mydb 127.0.0.1(60026)
182679 postgres  20   0 37.7g  23g  22g  S  15.3   7.4
2:11.07 postgres: pguser mydb 127.0.0.1(59112)
      58123 postgres  20   0 37.1g  36g  36g  R  13.6  11.7
13623:27 postgres: startup process
164415 postgres  20   0 37.7g  26g  26g  S  12.2   8.6
2:51.10 postgres: pguser mydb 127.0.0.1(42440)
      10674 root      20   0      0      0      0  S  11.7   0.0
22882:14 [shn_comp_wqa]
164421 postgres  20   0 37.7g  29g  28g  S  11.6   9.2
3:55.88 postgres: pguser mydb 127.0.0.1(42452)
      57570 root      20   0      0      0      0  S   6.0   0.0
6386:38 [flush-252:0]
      73195 postgres  20   0 37.1g 3532 2280  S   6.0   0.0
6207:28 postgres: wal receiver process
      58145 postgres  20   0 37.1g  36g  36g  S   3.6  11.6
4936:07 postgres: writer process
148192 postgres  20   0 37.7g  32g  31g  S   3.4  10.3
5:57.01 postgres: pguser mydb 127.0.0.1(53174)
      10683 root      20   0      0      0      0  S   2.1   0.0
3233:57 [shn_handle_luna]
164413 postgres  20   0 37.7g  28g  27g  S   2.1   9.0
3:31.48 postgres: pguser mydb 127.0.0.1(42436)
```

```
10681 root      20   0   0   0   0 S   1.7   0.0
1579:57 [shn_gc_wqa]
10675 root      20   0   0   0   0 S   1.1   0.0
886:09.06 [shn_wqa]
73174 postgres  20   0 184m 5816 1036 S   1.0   0.0
745:02.12 postgres: stats collector process
58144 postgres  20   0 37.2g 36g 36g S   0.3 11.6
812:08.97 postgres: checkpointer process
...
...
...
```

top

```
top - 12:01:03 up 93 days, 23:30,  1 user,  load average:
1.08, 1.09, 1.08
Tasks: 1042 total,   3 running, 1039 sleeping,   0 stopped,
0 zombie
Cpu(s):  5.3%us,  1.6%sy,  0.0%ni, 92.5%id,  0.2%wa,
0.0%hi,  0.3%si,  0.0%st
Mem:  330604220k total, 323235920k used,  7368300k free,
248996k buffers
Swap: 67108860k total,          0k used, 67108860k free,
295053464k cached
```

```
top - 12:01:03 932330,1
1515: 1.08, 1.09, 1.08
:1042,3, 1039,0,0
CPU:CPU5.3%,CPU1.6%,CPU0.0%,CPU92.5%,
IOCPU0.2%,CPU0.0%,CPU0.3%,0.0%
```

ss: 330604220k, sss323235920k, ss7368300k, buffers 248996k
ssss: 67108860k, sss0k, ss67108860k, cache 295053464k

ssssssssssssssssssssssssssssssss

·PID sssid

·USER sssss

·PR sssss

·NI sssssss

·VIRT sssssssssss

·RES sssssssssss

·SHR sssss

·S sssss D= sssss R= ss S= ss T= ss
ss/ss Z= sss

·%CPU sssss CPU sssss

·%MEM sssssssss

·TIME+ sssss CPU sssss 1/100

·COMMAND

top
PID USER PR NI VIRT
RES SHR S %CPU %MEM TIME+
COMMAND
top man top

2.free

free

[root@pghost1 ~]# free -g				
	total	used	free	shared
Mem:	315	306	9	36
0	278			
-/+ buffers/cache:		27	287	
Swap:	63	0	63	

Mem 315GB
306GB 9GB total=used+free
36GB buffers 0GB
244MB cache 278GB
buffers cache I/O
buffers cache
-/+ buffers/cache used-
buffers/cached

27GB free+buffers/cached
287GB
Swap Swap

free
Linux I/O
buffer cache page cache
buffer cache page cache
cache
Swap
Swap cache

```
[root@pghost1 ~]# sync  
[root@pghost1 ~]# echo 1 > /proc/sys/vm/drop_caches
```

3. vmstat

vmstat Linux
CPU vmstat

vmstat delay count delay count count delay
CTRL+C

```
[root@pghost1 ~]# vmstat 3 5
procs -----memory----- ---swap-- -----io----- --
system-- -----cpu-----
   r  b swpd  free   buff  cache  si  so    bi    bo    in
cs us sy id wa st
   3  0    0 7598968 243376 292787104    0    0  1354   984
0   0  4  2 94  1  0
   1  0    0 7491112 243388 292891072    0    0 27019 45264
45526 72922  4  2 94  0  0
   0  0    0 7434112 243408 292946016    0    0 11299 32468
45675 68756  4  1 94  0  0
   2  0    0 7348156 243468 293028032    0    0 19383 67697
48760 75728  4  2 94  0  0
   2  0    0 7289800 243480 293085504    0    0 12467 28441
46016 64854  4  1 95  0  0
```

vmstat

-r CPU 1

-b

-swpd

-free

-buff 缓冲区 buffer 缓冲区 buffer cache
缓冲区 缓冲区 cache

-cache 缓存 page cache 缓存
page cache

-si/so 交换 Swap 交换 Swap
交换

-bi/bo 块/字节 blocks/s

-in 输入

-cs CPU 核心数

4.iostat

iostat 工具 tty 设备 CD-
ROM 设备 iostat 工具
vmstat 工具 delay count
工具

```
[root@pghost1 ~]# iostat -dx /dev/dfa 5 5
Linux 2.6.32-696.el6.x86_64 (pghost1)      01/25/2018
_x86_64_      (48 CPU)
Device:      rrqm/s   wrqm/s     r/s     w/s   rsec/s
wsec/s avgrq-sz avgqu-sz   await r_await w_await  svctm
%util
dfa          0.00      0.00 3679.66 11761.85 129603.44
94094.80     14.49      0.17
```

	0.03	0.01	0.04	0.02	27.97		
Device:		rrqm/s	wrqm/s	r/s	w/s	rsec/s	wsec/s
avgrq-sz		avgqu-sz	await	r_await	w_await	svctm	%util
dfa	0.00	0.00	807.60	9619.20	33854.40	76953.60	10.63
	3.89						
	0.37	0.15	0.39	0.01	14.00		
Device:		rrqm/s	wrqm/s	r/s	w/s	rsec/s	wsec/s
avgrq-sz		avgqu-sz	await	r_await	w_await	svctm	%util
dfa	0.00	0.00	727.80	12904.40	29558.40	103235.20	9.74
	7.39						
	0.54	0.18	0.56	0.01	15.18		
Device:		rrqm/s	wrqm/s	r/s	w/s	rsec/s	wsec/s
avgrq-sz		avgqu-sz	await	r_await	w_await	svctm	%util
dfa	0.00	0.00	1682.00	13761.20	52283.20	110089.60	10.51
	16.64						
	1.07	0.15	1.18	0.01	20.76		
Device:		rrqm/s	wrqm/s	r/s	w/s	rsec/s	wsec/s
avgrq-sz		avgqu-sz	await	r_await	w_await	svctm	%util
dfa	0.00	0.00	609.80	23131.60	24737.60	185052.80	8.84
	16.50						
	0.69	0.24	0.70	0.01	21.62		

□□□□□□□□□□□□

-rrqm/s□wrqm/s□□□□□□□□□□□□□□OS□□□□□
□□□□□□□□□□□□

-r/s□w/s□□□□□□□□□□

-rsec/s□wsec/s□□□□□□□□□□□□

-avgrq-sz□□□□□□□□□□□□

-avgqu-sz 显示平均队列长度

-await 显示等待IO的时间

-svctm 显示服务时间

-%util 显示CPU利用率

5.mpstat

mpstat 显示CPU使用率

```
[root@pghost1 ~]# mpstat 5 5
Linux 2.6.32-696.el6.x86_64 (pghost1)      01/25/2018
_x86_64_      (48 CPU)
02:54:30 PM  CPU      %usr   %nice    %sys %iowait    %irq
%soft  %steal  %guest   %idle
02:54:35 PM  all       3.88    0.00    1.16    0.33    0.00
0.22    0.00    0.00   94.41
02:54:40 PM  all       4.06    0.00    0.98    0.14    0.00
0.24    0.00    0.00   94.58
02:54:45 PM  all       4.06    0.00    0.98    0.13    0.00
0.23    0.00    0.00   94.60
02:54:50 PM  all       4.15    0.00    1.37    0.37    0.00
0.24    0.00    0.00   93.88
02:54:55 PM  all       6.15    0.00    1.16    0.15    0.00
0.24    0.00    0.00   92.30
Average:      all       4.46    0.00    1.13    0.22    0.00
0.23    0.00    0.00   93.95
```

mpstat 显示CPU使用率
mpstat 显示CPU使用率 CPU 使用率 CPU 使用率

mpstat -P n 0 5 3
core 0 5
3

```
[root@pghost1 ~]# mpstat -P 0 5 3
Linux 2.6.32-696.el6.x86_64 (pghost1)          01/25/2018
_x86_64_ (48 CPU)
02:57:16 PM  CPU      %usr   %nice    %sys %iowait    %irq
%soft  %steal  %guest   %idle
02:57:21 PM    0   17.54    0.00    2.02    0.81    0.00
0.40    0.00    0.00   79.23
02:57:26 PM    0   19.07    0.00    2.43    0.61    0.00
0.41    0.00    0.00   77.48
02:57:31 PM    0   18.96    0.00    3.39    1.00    0.00
0.60    0.00    0.00   76.05
Average:        0   18.52    0.00    2.62    0.81    0.00
0.47    0.00    0.00   77.58
```

mpstat

-%usr

-%nice

-%sys

-%iowait

-%irq

-%soft

-%steal %guest

-%idle

6.sar

sar sar
/etc/cron.d/sysstat 10

```
[root@pghost1 ~]# cat /etc/cron.d/sysstat
# Run system activity accounting tool every 10 minutes
*/10 * * * * root /usr/lib64/sa/sa1 1 1
```

10
1

sar
AM/PM sar
LANG=C 24

sar

1 CPU

CPU

```
[root@pghost1 ~]# sar -q
12:00:01 AM      runq-sz    plist-sz      ldavg-1      ldavg-5      ldavg-
15
12:10:01 AM          9        1306         1.17         1.17
1.19
12:20:01 AM          4        1307         1.21         1.19
1.19
...
...
...
10:20:01 AM          9        1297         1.02         1.09
1.13
10:30:01 AM          7        1300         0.95         0.89
0.99
10:40:01 AM          8        1298         1.48         1.27
1.10
...
...
...
```

□□□□□□□□

-runq-sz□□□□□□□□

-plist-sz□□□□□□

-ldavg-1□ldavg-5□ldavg-15□□□□□5□□□□
15□□□□□□□□

□2□□□I/O□□

□□I/O□□□□□□□□□□

```
[root@pghost1 ~]# sar -b
12:00:01 AM          tps          rtps          wtps      bread/s
bwrtn/s
12:10:01 AM  13995.40      947.73   13047.68   35553.77
104404.64
12:20:01 AM  14389.13     1162.41   13226.72   40502.67
105837.12
12:30:01 AM  16420.53     1107.85   15312.69   43851.18
122524.99
...
...
...
01:00:01 PM  15961.31     1444.19   14517.12   52564.87
116160.35
01:10:01 PM  15327.42     1201.94   14125.48   43332.46
113027.64
...
...
...
```

□□□□□□□□□□

-tps□rtps□wtps□TPS□□

-bread/s□bwrtn/s□□□□□block□□□□□□□□□□
 □□□□block□□□□512□□□□□□□□□□□Block□□□□

□3□□□□□□□□□□

sar□□□□□□□□/var/log/sa/□□□□□□□□sar□□□□
 □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
 □/etc/sysconfig/sysstat□□HISTORY□□□□□□□□□□
 □□□□□28□□□□□/var/log/sa/□□□□□□□□□□□□□□□□

15:22
00:00:23:00:00 CPU

```
[root@pghost1 ~]# sar -q -f /var/log/sa/sa15 -s 22:00:00 -e 23:00:00
10:00:01 PM      runq-sz    plist-sz    ldavg-1    ldavg-5    ldavg-15
10:10:01 PM          5        1311        1.38        1.44
1.36
10:20:01 PM          2        1312        1.15        1.18
1.25
10:30:01 PM          5        1313        1.50        1.20
1.19
10:40:01 PM          4        1312        1.01        1.20
1.16
10:50:01 PM          6        1334        1.64        1.25
1.18
Average:          4        1316        1.34        1.25
1.23
```

sarLinux
gunplot
grafana

7.

nmontop
I/O CPU
OFFICE

iotop top 查看系统 I/O 使用情况
查看系统 I/O 使用情况 I/O 使用情况
查看系统 I/O 使用情况 pidstat 查看系统 I/O 使用情况

查看 iotop 查看系统 I/O 使用情况
查看系统 I/O 使用情况

查看系统 I/O 使用情况

-D 查看系统 I/O 使用情况

-R 查看系统 I/O 使用情况

-S 查看系统 I/O 使用情况

-T 查看系统 I/O 使用情况

-X dead 查看系统 I/O 使用情况

-Z 查看系统 I/O 使用情况

查看 D 查看系统 I/O 使用情况 “查看系统 I/O 使用情况”
查看系统 I/O 使用情况

```
[root@pghost1 ~]# for x in `seq 3`; do ps -eo state,pid,cmd  
| grep "^D"; echo "----"; sleep 5; done;
```

□

```
[root@pghost1 ~]# while true; do date; ps auxf | awk  
'{if($8=="D") print $0;}'; sleep 1; done;
```

□□□□pidstat-d 1□□□□□□□□□□□□□□□□I/O□
□□

10.2.2 Linux□□□I/O□□□□

□□I/O□□□□□□□□□□□□□□□□□□□□□□□□□□□□
Linux□□□□I/O□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□Linux□□□□□□□□□□
□□□Linux□□□□□□□□□□□□□□□□□□□□□□I/O□□□□□□□□□□
□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□

```
[root@pghost1 ~]# dmesg | grep -i scheduler  
io scheduler noop registered  
io scheduler anticipatory registered  
io scheduler deadline registered  
io scheduler cfq registered (default)
```

cfq□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□I/O□□□□□□□□□□□□□□□□□□□□□□□□I/O□□□□□□□□□□□□
□□□□Linux□□□□□□□□□□□□□□□□□□□□□□□□□□□□

noop 最简单的 I/O 调度器，采用 FIFO 调度策略，适用于 SSD。

deadline 调度器采用 FIFO 调度策略，适用于 SSD。它会在每个 I/O 请求到达时，为其分配一个 deadline，如果请求在 deadline 之前到达，则会被优先处理。

将 sda 的 I/O 调度器设置为 noop。

```
[root@pghost1 ~]# cat /sys/block/sda/queue/scheduler
noop anticipatory deadline [cfq]
```

将 sda 的 I/O 调度器设置为 noop。

将 shell 的 I/O 调度器设置为 noop。

```
[root@pghost1 ~]# echo noop > /sys/block/sda/queue/scheduler
[root@pghost1 ~]# cat /sys/block/sda/queue/scheduler
[noop] anticipatory deadline cfq
```

将 shell 的 I/O 调度器设置为 noop。

在 /etc/grub.conf 中添加以下配置：

10.2.3 配置

Linux I/O Linux I/O Linux I/O

```
[root@pghost1 ~]# /sbin/blockdev --getra /dev/sda
256
```

256 16384

```
[root@pghost1 ~]# /sbin/blockdev --setra 16384 /dev/sda
```

```
[root@pghost1 ~]# echo 16384
/sys/block/sda/queue/read_ahead_kb
```

/etc/rc.local

```
[root@pghost1 ~]# echo "/sbin/blockdev --setra 16384
/dev/dfa /dev/sda1" >> /etc/rc.local
```

10.2.4 五五五五

1.Swap

Swap free

```
[postgres@pghost1 ~]$ free -g
```

	total	used	free	shared
buffers	cached			
Mem:	378	34	344	30
0	30			
-/+ buffers/cache:		3	375	
Swap:	31	0	31	

0000000000000000Swap000031GB0000
0000000000000000Swap0

00000vmstat0000000000000000

```
[postgres@pghost1 ~]$ vmstat 5 3
procs -----memory----- ---swap-- -----io----- --
system-- -----cpu-----
      r  b   swpd   free   buff   cache   si    so    bi    bo
in   cs us sy id wa st
    0   0     0 361127872 390616 31882912    0    0     2   11
0     0   0   0 100   0   0
    0   0     0 361127840 390620 31882912    0    0     0   30
744  533   0   0 100   0   0
    1   0     0 361127840 390620 31882912    0    0     0   22
1161  694   0   0 100   0   0
[postgres@pghost1 ~]$
```

0vmstat000000000swpd0000000000
0Swap0000000000000000Swap0

0000000000000000Swap0000000000000000
0000Swap0000Swap0000000000000000Swap0000
Swap00000000Swap0000swapoff0000Swap00
00swapon00000

```
[root@pghost1 ~]# swapoff -a
[root@pghost1 ~]#
[root@pghost1 ~]# free | grep Swap
Swap:          0          0          0
[root@pghost1 ~]#
[root@pghost1 ~]# swapon -a
[root@pghost1 ~]#
```

```
[root@pghost1 ~]# free | grep Swap
Swap:      33554428          0    33554428
```

Swap空间使用0 Swap
Swap空间

2. 配置

Transparent HugePages配置
配置Transparent HugePages
配置Transparent HugePages

配置Transparent HugePages

```
[root@pghost1 ~]# cat
/sys/kernel/mm/transparent_hugepage/enabled
[always] madvise never
```

配置Transparent HugePages
配置Transparent HugePages

```
[root@pghost1 ~]# echo never >
/sys/kernel/mm/transparent_hugepage/enabled
[root@pghost1 ~]# cat
/sys/kernel/mm/transparent_hugepage/enabled
always madvise [never]
```

配置Transparent HugePages/etc/rc.local配置

```
if test -f /sys/kernel/mm/transparent_hugepage/enabled; then
    echo never > /sys/kernel/mm/transparent_hugepage/enabled
fi
if test -f /sys/kernel/mm/transparent_hugepage/defrag; then
    echo never > /sys/kernel/mm/transparent_hugepage/defrag
fi
```

```
kernel /boot/vmlinuz-2.6.32-642.11.1.el6.x86_64 ro
transparent_hugepage=never
```

```
kernel /boot/vmlinuz-2.6.32-642.11.1.el6.x86_64 ro
root=UUID=c429e8ae-c35d-4bc0-a781-17bbb95a75cf nomodeset
rd_NO_LUKS KEYBOARDTYPE=pc KEYTABLE=us LANG=en_US.UTF-8
rd_NO_MD SYSFONT=latarcyrheb-sun16 crashkernel=auto
rd_NO_LVM rd_NO_DM rhgb quiet numa=off
transparent_hugepage=never
```

3.NUMA

```
NUMA local CPU local
local local local
CPU NUMA
NUMA
```

```
NUMA
```

```
[root@pghost1 ~]# numactl --hardware
available: 2 nodes (0-1)
node 0 cpus: 0 2 4 6 8 10 12 14 16 18 20 22 24 26 28 30 32
```



```
34 36 38 40 42 44 46
node 0 size: 196514 MB
node 0 free: 186290 MB
node 1 cpus: 1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 31 33
35 37 39 41 43 45 47
node 1 size: 196608 MB
node 1 free: 192336 MB
node distances:
node      0      1
  0:    10    21
  1:    21    10
```

numastat

```
[root@pghost1 ~]# numastat
```

	node0	node1
numa_hit	27207118	27526494
numa_miss	0	0
numa_foreign	0	0
interleave_hit	111148	111125
local_node	27205425	27405472
other_node	1693	121022

numactl CPU

available 2nodes 0-1

NUMA BIOS

Node Interleaving

Node Interleaving

Enabled

NUMA Disabled
NUMA Disabled
NUMA

kernel /etc/grub.conf
numa=off NUMA

```
kernel /boot/vmlinuz-2.6.32-642.11.1.el6.x86_64 ro
root=UUID=c429e8ae-c35d-4bc0-a781-17bbb95a75cf nomodeset
rd_NO_LUKS KEYBOARDTYPE=pc KEYTABLE=us LANG=en_US.UTF-8
rd_NO_MD SYSFONT=latarcyrheb-sun16 crashkernel=auto
rd_NO_LVM rd_NO_DM rhgb quiet numa=off
```

NUMA

```
[root@pghost1 ~]# numactl --hardware
available: 1 nodes (0)
node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38
39 40 41 42 43 44 45 46 47
node 0 size: 393122 MB
node 0 free: 379902 MB
node distances:
node 0
  0: 10
[root@pghost1 ~]#
[root@pghost1 ~]# numastat
```

	node0
numa_hit	3613403
numa_miss	0
numa_foreign	0
interleave_hit	222419
local_node	3613403
other_node	0

A diagram consisting of a 4x28 grid of squares. The word "Linux" is centered in the second row. The top row has 28 squares, the second row has 28 squares with "Linux" in the center, the third row has 28 squares, and the fourth row has 3 squares.

10.3 数据库

10.3.1 数据库

postgresql.conf 数据库配置参数文件
数据库配置参数文件位于数据库安装目录的
data目录下

PostgreSQL 数据库配置参数文件
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数
shared_buffers 数据库配置参数

postgresql.conf 数据库配置参数文件
PostgreSQL 10 数据库配置参数
128MB 数据库配置参数
PostgreSQL 数据库配置参数
shared_buffers 数据库配置参数
pgbench 数据库配置参数

work_mem 参数控制 hash 排序时使用的内存大小。默认值为 4MB。可以通过修改 work_mem 参数来调整 hash 排序的内存使用量。PostgreSQL 支持 Top-N heapsort、Quick sort、External merge 和 External merge 排序。work_mem 参数

random_page_cost 参数控制随机访问页面的成本。默认值为 4。可以通过修改 random_page_cost 参数来调整随机访问页面的成本。seq_page_cost 参数控制顺序访问页面的成本。默认值为 1.5。可以通过修改 seq_page_cost 参数来调整顺序访问页面的成本。

PostgreSQL 支持多种排序算法。可以通过修改相关参数来调整排序算法的性能。

10.3.2 索引

PostgreSQL 支持多种索引。可以通过修改相关参数来调整索引的性能。pg_stat 和 pg_statio 参数用于监控数据库的性能。

autovacuum
pg_statistics
pg_stats

1.pg_stat_database

pg_stat_database

```
mydb=# \d pg_stat_database
View "pg_catalog.pg_stat_database"
      Column      |          Type          | Collation |
Nullable | Default |
-----+-----+-----+-----
datid      | oid      |           |
|
datname    | name     |           |
|
numbackends | integer  |           |
|
xact_commit | bigint   |           |
|
xact_rollback | bigint   |           |
|
blks_read   | bigint   |           |
|
blks_hit    | bigint   |           |
|
tup_returned | bigint   |           |
|
tup_fetched | bigint   |           |
|
tup_inserted | bigint   |           |
|
```

tup_updated	bigint		
tup_deleted	bigint		
conflicts	bigint		
temp_files	bigint		
temp_bytes	bigint		
deadlocks	bigint		
blk_read_time	double precision		
blk_write_time	double precision		
stats_reset	timestamp with time zone		

□□□□□□□

·numbackends□□□□□□□□□□□□□□□□□□
cpu□□□1.5□□□□□□□□□□□□

·blks_read□blks_hit□□□□□□□□□□□□□□□□
□□□□□

·xact_commit□xact_rollback□□□□□□□□□□
□□

·deadlocks□□□□□□pg_stat_reset□□□□□□□
□□

□□□□□□□□□□□□□□□□□□

```
SELECT blks_hit::float/(blks_read + blks_hit) as  
cache_hit_ratio FROM pg_stat_database WHERE  
datname=current_database();
```

IO 1
shared_buffers 99%

```
SELECT xact_commit::float/(xact_commit + xact_rollback) as  
successful_xact_ratio FROM pg_stat_database WHERE  
datname=current_database();
```

1

pg_stat_database
numbackends stats_reset
stats_reset pg_stat_reset
pg_stat_reset
“stat” “reset”
pg_statistics
ANALYZE pg_stat_reset

2.pg_stat_user_tables

pg_stat_user_tables
pg_stat_user_tables

```
mydb=# \d pg_stat_user_tables
```

```
View "pg_catalog.pg_stat_user_tables"
Column | Type | Collation |
Nullability | Default |
-----+-----+-----+
relid | oid | |
|
schemaname | name | |
|
relname | name | |
|
seq_scan | bigint | |
|
seq_tup_read | bigint | |
|
idx_scan | bigint | |
|
idx_tup_fetch | bigint | |
|
n_tup_ins | bigint | |
|
n_tup_upd | bigint | |
|
n_tup_del | bigint | |
|
n_tup_hot_upd | bigint | |
|
n_live_tup | bigint | |
|
n_dead_tup | bigint | |
|
n_mod_since_analyze | bigint | |
|
last_vacuum | timestamp with time zone | |
|
```

last_autovacuum	timestamp with time zone	
last_analyze	timestamp with time zone	
last_autoanalyze	timestamp with time zone	
vacuum_count	bigint	
autovacuum_count	bigint	
analyze_count	bigint	
autoanalyze_count	bigint	

last_vacuum last_analyze
vacuum analyze

last_autovacuum last_autoanalyze
autovacuum autovacuum
analyze

idx_scan idx_tup_fetch
idx_tup_fetch

seq_scan seq_tup_read
seq_tup_read

n_tup_ins n_tup_upd n_tup_del
n_tup_ins

```

n_live_tup n_dead_tup live tuple
dead tuple

```

VS

```
SELECT sum(idx_scan)/(sum(idx_scan) + sum(seq_scan)) as
idx_scan_ratio FROM pg_stat_all_tables WHERE
schemaname='your_schema';
SELECT relname,idx_scan::float/(idx_scan+seq_scan+1) as
idx_scan_ratio FROM pg_stat_all_tables WHERE
schemaname='your_schema' ORDER BY idx_scan_ratio ASC;
```

[illegible]

3.pg_stat_statements

```
pg_stat_statements
postgres auto explain
```

```
pg_stat_statements
postgresql.conf
```

```
shared_preload_libraries = 'pg_stat_statements'
pg_stat_statements.track = all
```

CREATE EXTENSION

```
mydb=# CREATE EXTENSION pg_stat_statements;
CREATE EXTENSION
```

pg_stat_statements

```
mydb=# \d pg_stat_statements
          View "public.pg_stat_statements"
      Column |          Type          | Collation |
Nullable | Default |
-----+-----+-----+-----
--+-+-----
userid      | oid      |          |
|
dbid        | oid      |          |
|
queryid     | bigint   |          |
|
query       | text     |          |
|
calls       | bigint   |          |
|
total_time  | double precision |          |
|
min_time    | double precision |          |
|
max_time    | double precision |          |
|
mean_time   | double precision |          |
|
stddev_time | double precision |          |
|
rows        | bigint   |          |
```



```
126 | 365.336761904762 | SELECT tableoid, oid,  
proname, prolang, pronargs, proargtypes, prorettytype, proac  
(3 rows)
```

```
pg_stat_statements  
pg_stat_statements_reset  
pg_stat_statements
```

4. SQL

```
PostgreSQL EXPLAIN
```

```
mydb=# EXPLAIN SELECT * FROM tbl;  
          QUERY PLAN
```

```
-----  
Seq Scan on tbl  (cost=0.00..20.70 rows=1070 width=48)  
(1 row)
```

```
EXPLAIN ANALYZE
```

```
mydb=# EXPLAIN ANALYZE SELECT * FROM tbl;  
          QUERY PLAN
```

```
-----  
Seq Scan on tbl  (cost=0.00..20.70 rows=1070 width=48)  
(actual time=0.021..0.023 rows=3 loops=1)
```

Planning time: 0.117 ms
Execution time: 0.058 ms
(3 rows)

ANALYZE
INSERT
UPDATE
DELETE
CREATE TABLE
AS
EXECUTE
ANALYZE

```
mydb=# BEGIN;
BEGIN
mydb=# EXPLAIN ANALYZE UPDATE tbl SET ival = ival * 10 WHERE
id = 1;
```

QUERY PLAN

```
-----
Update on tbl  (cost=0.15..8.17 rows=1 width=54) (actual
time=0.159..0.159 rows=0 loops=1)
  -> Index Scan using tbl_pkey on tbl  (cost=0.15..8.17
rows=1 width=54) (actual time=0.046..0.047 rows=1 loops=1)
        Index Cond: (id = 1)
Planning time: 4.237 ms
Execution time: 0.315 ms
(5 rows)
mydb=# ROLLBACK;
ROLLBACK
```

Execution time
0.315ms
Planning time
4.237ms
tbl

tbl_pkey Index Scan
Update
cost=0.00..xxx
rows=xxx
PostgreSQL
width=1917
ANALYZE

ANALYZE COSTS
BUFFERS TIMING FORMAT

```
mydb=# EXPLAIN (ANALYZE on, TIMING on , VERBOSE on, BUFFERS on) SELECT * FROM tbl WHERE id = 10;
               QUERY PLAN
```

```
-----
Index Scan using tbl_pkey on public.tbl (cost=0.15..8.17
rows=1 width=48) (actual time=0.015..0.015 rows=0 loops=1)
  Output: id, ival, description, created_time
  Index Cond: (tbl.id = 10)
  Buffers: shared hit=1
Planning time: 0.177 ms
Execution time: 0.065 ms
(6 rows)
```

EXPLAIN ANALYZE COSTS
BUFFERS TIMING VERBOSE

LOG: PARSE ANALYSIS STATISTICS

DETAIL: ! system usage stats:

! 0.000086 elapsed 0.000000 user 0.000000 system sec
! [0.061990 user 0.014997 sys total]
! 0/0 [600/0] filesystem blocks in/out
! 0/0 [0/2640] page faults/reclaims, 0 [0] swaps
! 0 [0] signals rcvd, 0/0 [0/0] messages rcvd/sent
! 0/0 [55/1] voluntary/involuntary context switches

LOG: REWRITER STATISTICS

DETAIL: ! system usage stats:

! 0.000001 elapsed 0.000000 user 0.000000 system sec
! [0.061990 user 0.014997 sys total]
! 0/0 [600/0] filesystem blocks in/out
! 0/0 [0/2640] page faults/reclaims, 0 [0] swaps
! 0 [0] signals rcvd, 0/0 [0/0] messages rcvd/sent
! 0/0 [55/1] voluntary/involuntary context switches

LOG: PLANNER STATISTICS

DETAIL: ! system usage stats:

! 0.000165 elapsed 0.001000 user 0.000000 system sec
! [0.063990 user 0.014997 sys total]
! 0/0 [600/0] filesystem blocks in/out
! 0/0 [0/2640] page faults/reclaims, 0 [0] swaps
! 0 [0] signals rcvd, 0/0 [0/0] messages rcvd/sent
! 0/0 [58/1] voluntary/involuntary context switches

QUERY PLAN

Limit (cost=0.00..0.43 rows=10 width=224) (actual
time=0.028..0.034 rows=10 loops=1)
-> Seq Scan on tbl (cost=0.00..64771378.60
rows=1496764160 width=224) (actual time=0.025..0.027 rows=10
loops=1)
Planning time: 0.198 ms
Execution time: 0.083 ms
(4 rows)

□□□□□□□□□□□□□□□□

mydb=# set client_min_messages = log;
mydb=# set log_parser_stats = off;
mydb=# set log_planner_stats = off;
mydb=# set log_statement_stats = on;

```

mydb=# EXPLAIN ANALYZE select * from tbl limit 10;
LOG:  QUERY STATISTICS
DETAIL:  ! system usage stats:
!       0.000603 elapsed 0.000000 user 0.000000 system sec
!       [0.065989 user 0.014997 sys total]
!       0/0 [600/0] filesystem blocks in/out
!       0/0 [0/2640] page faults/reclaims, 0 [0] swaps
!       0 [0] signals rcvd, 0/0 [0/0] messages rcvd/sent
!       0/0 [62/1] voluntary/involuntary context switches
          QUERY PLAN

-----
Limit  (cost=0.00..0.43 rows=10 width=224) (actual
time=0.027..0.033 rows=10 loops=1)
  -> Seq Scan on tbl  (cost=0.00..64771378.60
rows=1496764160 width=224) (actual time=0.026..0.026 rows=10
loops=1)
    Planning time: 0.154 ms
    Execution time: 0.082 ms
(4 rows)

```

```

    parser planner
    DETAIL system usage stats
    CPU 3
    I/O 4
    5 IPC 6

```

10.3.3 索引

```

    PostgreSQL B-tree
    Hash GiST SP-GiST GIN BRIN Bloom
    B-tree B-tree CREATE
    INDEX B-tree

```

HashGIN
GiST
GiSTB-treeR-tree
PostGISGiST
PostgreSQLB-treeGiSTGINBRIN32

PostgreSQLCREATE INDEX
CONCURRENTLY
CONCURRENTLY
PostgreSQLMVCC
“”CREATE
INDEX CONCURRENTLY
CREATE INDEX

```
mydb=# CREATE UNIQUE INDEX CONCURRENTLY ON mytbl USING  
btree(id);  
CREATE INDEX
```

mytbl_id_idx mytbl_pkey

```
mydb=# SELECT
schemaname,relname,indexrelname,pg_relation_size(indexrelid)
AS index_size,idx_scan,idx_tup_read,idx_tup_fetch FROM
pg_stat_user_indexes WHERE indexrelname IN (SELECT indexname
FROM pg_indexes WHERE schemaname = 'public' AND tablename =
'mytbl');
schemaname | relname | indexrelname | index_size | idx_scan
| idx_tup_read | idx_tup_fetch
-----
-
 public    | mytbl  | mytbl_pkey   | 223051776 | 1403532 |
1413850 |      1403532
 public    | mytbl  | mytbl_id_idx | 222887936 |          0 |
0 |          0
(2 rows)
```

```
mydb=# BEGIN;
BEGIN
mydb=# ALTER TABLE mytbl DROP CONSTRAINT mytbl_pkey;
ALTER TABLE
mydb=# ALTER TABLE mytbl ADD CONSTRAINT mytbl_id_idx PRIMARY
KEY USING INDEX mytbl_id_idx;
ALTER TABLE
mydb=# END;
COMMIT
```


10.4 □□□□

```

#####
Linux#####I/O#####
#####Swap#####NUMA#####
#####
#####
#####
#####
#####VACUUM#####ANALYZE#####
VACUUM FREEZE#####
#####
#####
#####
#####
#####

```

11 pgbench

pgbench is a tool for benchmarking PostgreSQL database systems.
It is a client-side benchmark that simulates a workload of transactions.
The transactions are defined in a script file, and the benchmark can be run
against a PostgreSQL database. The benchmark can be used to measure
the performance of a database system under various workloads.
The benchmark can be run against a PostgreSQL database, or it can be
run against a PostgreSQL database that is running on a remote host.
The benchmark can be run against a PostgreSQL database that is running
on a remote host, or it can be run against a PostgreSQL database that is
running on a remote host. The benchmark can be run against a PostgreSQL
pgbench

11.1

- CPU가 데이터를 처리하는 방식
- 데이터의 저장과 검색 방법
- 데이터의 전송과 수신 방법
- 데이터의 암호화와 복호화 방법

PostgreSQL

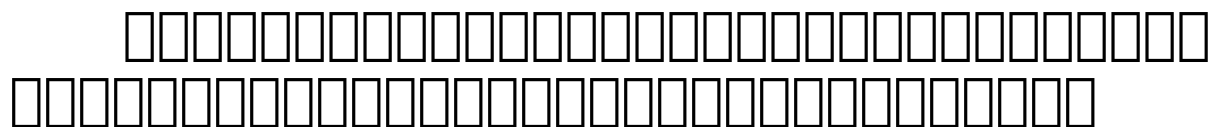
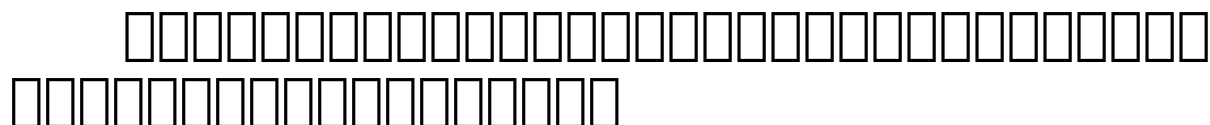
11.1.1 11.1.1.1 11.1.1.1.1 11.1.1.1.1.1

通过压测工具对系统性能进行测试，主要关注以下指标：
Throughput（吞吐量）、RT（响应时间）、Latency（延迟）等。

TPS（每秒事务数）是衡量系统性能的重要指标之一。在压测过程中，需要记录不同负载下的 TPS 值。通常，系统在高负载下会出现性能下降，例如 TPS 可能会从 90% 的峰值下降到 10% 左右，甚至出现 5 秒以上的响应延迟。通过对比不同配置下的性能表现，可以评估系统的稳定性和扩展性。

11.1.3 性能优化策略

性能优化的核心目标是提升系统的吞吐量和降低延迟。常见的优化策略包括：
1. 数据库优化：通过索引优化、SQL 语句调优等方式提升数据库性能。
2. 缓存策略：引入 Redis 等缓存系统，减少数据库压力。
3. 负载均衡：通过 Nginx 等工具实现流量均衡，避免单点过载。
4. 硬件升级：根据压测结果，适当增加 CPU、内存或 I/O 资源。
5. 代码优化：对热点代码进行重构，提升执行效率。
6. 异步处理：利用消息队列（如 RabbitMQ）处理异步任务，提升并发能力。
7. 限流熔断：防止系统因突发流量而崩溃，保障核心业务可用性。



11.2 pgbench

TPC Transaction Processing Performance Council
<http://www.tpc.org> TPC-A TPC-B
TPC-C TPC-D TPC-E TPC-W
TPC-C OLTP
OLTP TPC-E
benchmarks sql sysbench
PostgreSQL pgbench
pgbench TPC-B

11.2.1 pgbench

pgbench

```
transaction type: <builtin: TPC-B (sort of)>
scaling factor: 100
query mode: simple
number of clients: 1
number of threads: 1
number of transactions per client: 10
number of transactions actually processed: 10/10
latency average = 2.557 ms
tps = 391.152261 (including connections establishing)
tps = 399.368200 (excluding connections establishing)
```

transaction type

scaling factor pgbench

query mode simple
extended prepared

number of clients

number of threads

number of transactions per client

number of transactions actually
processed

latency average

TPS TPS

11.2.2

1.

```
pgbench[4]
pgbench_branches pgbench_tellers
pgbench_accounts pgbench_history
pgbench pgbench
[ ]
pgbench
```

pgbench

`pgbench -i [OPTION]... [DBNAME]`

pgbench

[illegible]

```
pgbench_accounts      100000k
pgbench_history        0
--foreign-keys         4000000000000000
--index-tablespace=TABLESPACE
                        0000000000
--tablespace=TABLESPACE 0000000000
--unlogged-tables      40000000UNLOGGED
```

0000000000000000

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -i -s 2 -F
80 -h pghost1 -p 1921 -U pguser -d mydb
creating tables...
100000 of 200000 tuples (50%) done (elapsed 0.02 s,
remaining 0.02 s)
200000 of 200000 tuples (100%) done (elapsed 0.16 s,
remaining 0.00 s)
vacuum...
set primary keys...
done.
```

2. pgbench 00000000

pgbench 00000000 tpcb-like simple-
update select-only 0000000000000000
pgbench 000000000000

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -b list
Available builtin scripts:
    tpcb-like
    simple-update
    select-only
```

tpcb-likeSELECT
UPDATEINSERT

```
BEGIN;
    UPDATE pgbench_accounts SET abalance = abalance + :delta
WHERE aid = :aid;
    SELECT abalance FROM pgbench_accounts WHERE aid = :aid;
    UPDATE pgbench_tellers SET tbalance = tbalance + :delta
WHERE tid = :tid;
    UPDATE pgbench_branches SET bbalance = bbalance + :delta
WHERE bid = :bid;
    INSERT INTO pgbench_history (tid, bid, aid, delta,
mtime) VALUES (:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);
END;
simple-update
BEGIN;
    UPDATE pgbench_accounts SET abalance = abalance + :delta
WHERE aid = :aid;
    SELECT abalance FROM pgbench_accounts WHERE aid = :aid;
    INSERT INTO pgbench_history (tid, bid, aid, delta,
mtime) VALUES (:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);
END;
select-only
BEGIN;
    SELECT abalance FROM pgbench_accounts WHERE aid = :aid;
END;
```

```
-b scriptname[@weight]
--builtin = scriptname[@weight]
```

scriptname
tpcb-likesimple-update

XXXXXXXXXX

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -b simple-  
update -h pghost1 -p 1921 -U pguser mydb  
starting vacuum...end.  
transaction type: <builtin: simple update>  
scaling factor: 100  
query mode: simple  
number of clients: 1  
number of threads: 1  
number of transactions per client: 10  
number of transactions actually processed: 10/10  
latency average = 1.631 ms  
tps = 613.238093 (including connections establishing)  
tps = 631.101768 (excluding connections establishing)
```

XXXXXXXX3XXXXXXXXXXXXXXXXXXXXX@
XX@XXXXXXXXXXXXXXXXXXXXXXXXXXXXsimple-
updateselect-onlyXXXXXXXX2X8XXXXXXXX
XXXXXXXXXXXX

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -b simple-  
update@2 -b select-only@8 -b tpcb@0 -h pghost1 -p 1921 -U  
pguser mydb  
starting vacuum...end.  
transaction type: multiple scripts  
scaling factor: 100  
query mode: simple  
number of clients: 1  
number of threads: 1  
number of transactions per client: 10  
number of transactions actually processed: 10/10  
latency average = 0.617 ms  
tps = 1621.779592 (including connections establishing)  
tps = 1753.156910 (excluding connections establishing)  
SQL script 1: <builtin: TPC-B (sort of)>  
- weight: 0 (targets 0.0% of total)
```

- 0 transactions (0.0% of total, tps = 0.000000)
- latency average = -nan ms
- latency stddev = -nan ms

SQL script 2: <builtin: select only>

- weight: 8 (targets 80.0% of total)
- 9 transactions (90.0% of total, tps = 1459.601633)
- latency average = 0.439 ms
- latency stddev = 0.137 ms

SQL script 3: <builtin: simple update>

- weight: 2 (targets 20.0% of total)
- 1 transactions (10.0% of total, tps = 162.177959)
- latency average = 1.738 ms
- latency stddev = 0.000 ms

```

@
tpcb-like
t simple select s
si se

```

ambiguous builtin name: 2 builtin scripts found for prefix "s"

```

-b simple-update-N--skip-some-
updates-b select-only-S--select-
only

```

```

TPS

```

11.2.3 pgbench

pgbenchはPostgreSQLのベンチマークツール。pgbenchは、データベースの性能を測定するためのツール。pgbenchは、データベースの性能を測定するためのツール。pgbenchは、データベースの性能を測定するためのツール。

1. テーブルの作成

テーブルの作成

```
CREATE TABLE tbl
(
    id SERIAL PRIMARY KEY,
    ival INT
);
```

2. pgbenchの実行

pgbenchの実行

```
[postgres@pghost2 ~]$ echo "SELECT id,ival FROM tbl ORDER BY
id DESC LIMIT 10;" > bench_script_for_select.sql
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -f
bench_script_for_select.sql -h pghost1 -p 1921 -U pguser
mydb
transaction type: bench_script_for_select.sql
...
tps = 3448.671865 (including connections establishing)
tps = 4097.559616 (excluding connections establishing)
```

SQLデータベースに接続し、
psqlコマンドを実行し、
pgbenchを実行する。

\sleep number[us|ms|s] 実行時間
を指定する。

\set varname expression 変数に
値を設定する。pgbenchでは
random int int を指定する。

tbl ival を指定する。
実行。

```
[postgres@pghost2 ~]$ cat bench_script_for_insert.sql
\sleep 500 ms
\set ival random(1, 100000)
INSERT INTO tbl(ival) VALUES (:ival);
```

pgbenchを実行する。

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -f
bench_script_for_insert.sql -h pghost1 -p 1921 -U pguser
mydb
starting vacuum...end.
transaction type: bench_script_for_insert.sql
...
latency average = 501.291 ms
number of transactions actually processed: 14846
...
```

pgbench

```
mydb=# SELECT COUNT(*) FROM tbl;
count
-----
 14846
(1 row)
mydb=# SELECT id,ival FROM tbl ORDER BY id DESC LIMIT 3;
 id   | ival
-----+-----
 14846 |  95018
 14845 |  88153
 14844 |  21896
(3 rows)
```

pgbench INSERT 14846 latency average 500 14846

pgbench -f @ 3 10

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -T 60 -f
bench_script_for_insert.sql@3 -f
bench_script_for_insert.sql@10 -h pghost1 -p 1921 -U pguser
mydb
starting vacuum...end.
transaction type: multiple scripts
scaling factor: 1
query mode: simple
number of clients: 1
number of threads: 1
```

```
duration: 60 s
number of transactions actually processed: 176523
latency average = 0.340 ms
tps = 2942.048332 (including connections establishing)
tps = 2942.071422 (excluding connections establishing)
SQL script 1: bench_script_for_insert.sql
- weight: 3 (targets 23.1% of total)
- 40921 transactions (23.2% of total, tps = 682.016280)
- latency average = 0.340 ms
- latency stddev = 0.035 ms
SQL script 2: bench_script_for_insert.sql
- weight: 10 (targets 76.9% of total)
- 135602 transactions (76.8% of total, tps =
2260.032052)
- latency average = 0.340 ms
- latency stddev = 0.038 ms
```

11.2.4 配置

1. 配置

```
-c 配置
1 -c 配置
配置
```

配置4配置

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -c 4 -h
pghost1 -p 1921 -U pguser mydb
...
number of clients: 4
...
latency average = 4.826 ms
tps = 828.781956 (including connections establishing)
tps = 894.930906 (excluding connections establishing)
```

pgbench

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -c 4 -C -h  
pghost1 -p 1921 -U pguser mydb
```

```
...  
number of clients: 4
```

```
...  
latency average = 19.627 ms  
tps = 203.797152 (including connections establishing)  
tps = 256.799857 (excluding connections establishing)
```

pgbench CPU usage -j
pgbench

2. pgbench

pgbench

·-T seconds --time=seconds
pgbench 1 --time=3600

·-t transactions --
transactions=transactions
pgbench 10

pgbench
pgbench

pgbench 简介

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -c 4 -h  
pghost1 -p 1921 -U pguser mydb
```

```
...  
number of clients: 4  
...  
number of transactions actually processed: 40/40  
...
```

pgbench 4 客户端同时运行，每个客户端执行 40 次事务。

3. 性能测试

pgbench -R 选项用于测试 TPS（每秒事务数）。

4. 延迟测试

pgbench -L 选项用于测试延迟（ms）。

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -T 10 -L 1 -  
c 8 -j 8 -f bench_script_for_select.sql@3 -f  
bench_script_for_update.sql@10 -h pghost1 -p 1921 -U pguser  
mydb
```

```
...  
number of transactions actually processed: 129234  
number of transactions above the 1.0 ms latency limit: 2226
```


pgbench_log
log-prefix prefix_name

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pgbench -T 10 -l --  
log-prefix=custom -c 6 -j 2 -f bench_script_for_insert.sql@3  
-f bench_script_for_insert.sql@10 -h pghost1 -p 1921 -U  
pguser mydb
```

pgbench
custom.151940 custom2.151940.2

client_id	transaction_no	time	script_no	time_epoch	time_us	[schedule_lag]
5	17	1294	1	1515153407	106478	
4	15	322	0	1515152900	291756	
3	7	357	0	1515152900	291768	
...						
4	16	334	0	1515152900	292090	
...						

client_id id transaction_no
time
script_no
time_epoch/time_us Unix
ISO
8601

11.3 □□□□

[illegible]

12 数据库复制

PostgreSQL 9.0 支持两种复制：
Streaming Replication（流复制）
和 Logical Replication（逻辑复制）。
pghost1 是 PostgreSQL 主数据库，
pgghost2 是 PostgreSQL 备数据库。
Primary Database Master pghost2
Standby Database Slave pghost1
pgghost2 是主数据库，pgghost1 是备数据库。
pgghost2 是主数据库，pgghost1 是备数据库。
pgghost2 是主数据库，pgghost1 是备数据库。

Logical Replication（逻辑复制）
是 PostgreSQL 10 引入的。
它使用 Slony-I 或 Londiste 或 pglogical 实现。
PostgreSQL 10 支持逻辑复制。
PostgreSQL 10 支持逻辑复制。
PostgreSQL 10 支持逻辑复制。

WAL（Write-Ahead Logging）
是 PostgreSQL 的预写式日志。

Diagram illustrating a memory layout with four rows of blocks. Each row contains 16 blocks. The first row has 'WAL' in the 15th block. The second row has 'WAL' in the 15th block. The third row has 'WAL' in the 1st and 15th blocks. The fourth row has 'WAL' in the 1st and 5th blocks.

```

      . WAL
WAL WAL WAL
WAL WAL WAL
WAL DML
WAL

```

· PostgreSQL

· **DDL** **DDL**

[illegible]

· PostgreSQL

PostgreSQL

Delayed Standbys Quorum Commit

12.1 部署环境

部署环境要求如下：
1. 操作系统：CentOS 6.9
2. 数据库：PostgreSQL 10
3. 网络：两台服务器通过内网连接，且能够访问互联网
4. 磁盘：每台服务器至少有一个 100GB 的磁盘空间
5. 用户：创建两个普通用户，分别用于主节点和备节点的部署

部署环境要求如下：
1. 操作系统：CentOS 6.9
2. 数据库：PostgreSQL 10
3. 网络：两台服务器通过内网连接，且能够访问互联网
4. 磁盘：每台服务器至少有一个 100GB 的磁盘空间
5. 用户：创建两个普通用户，分别用于主节点和备节点的部署

部署环境要求

12-1 部署环境要求

主机	主机名	IP 地址	操作系统	PostgreSQL 版本
主节点	pghost1	192.168.28.74	CentOS6.9	PostgreSQL10
备节点	pghost2	192.168.28.75	CentOS6.9	PostgreSQL10

部署环境要求如下：
1. 操作系统：CentOS 6.9
2. 数据库：PostgreSQL 10
3. 网络：两台服务器通过内网连接，且能够访问互联网
4. 磁盘：每台服务器至少有一个 100GB 的磁盘空间
5. 用户：创建两个普通用户，分别用于主节点和备节点的部署

12.1.1 部署环境要求

创建PostgreSQL用户1
PostgreSQL用户1的密码为pghost1
pghost2创建PostgreSQL

创建pghost1
pghost2用户

```
# groupadd postgres
# useradd postgres -g postgres
# passwd postgres
# mkdir -p /database/pg10/pg_root
# mkdir -p /database/pg10/pg_tbs
# chown -R postgres:postgres /database/pg10
```

/database/pg10/pg_root
/database/pg10/pg_tbs

创建postgres用户
postgres用户的主目录为/home/postgres/.bash_profile

```
export PGPOR=1921
export PGUSER=postgres
export PGDATA=/database/pg10/pg_root
export LANG=en_US.utf8
export PGHOME=/opt/pgsql
export
LD_LIBRARY_PATH=$PGHOME/lib:/lib64:/usr/lib64:/usr/local/lib
64:/lib:/usr/lib:/usr/local/lib
export PATH=$PGHOME/bin:$PATH:.
export MANPATH=$PGHOME/share/man:$MANPATH
alias rm='rm -i'
alias ll='ls -lh'
```

编译安装PostgreSQL

编译安装PostgreSQL

```
# tar jxvf postgresql-10.0.tar.bz2
# cd postgresql-10.0
# ./configure --prefix=/opt/pgsql_10.0 --with-pgport=1921
```

configure 编译选项 zlib readline
configure 编译选项 yum install
编译选项 zlib readline

```
# yum install zlib readline
```

编译安装PostgreSQL

```
# gmake world
# gmake install-world
```

gmake world
gmake
gmake install-world
PostgreSQL /opt/pgsql_10.0
/opt/pgsql_10.0/share/doc
/opt/pgsql_10.0/share/extension

PostgreSQL을
gmake world

PostgreSQL pgsql_10.0

```
# ln -s /opt/pgsql_10.0 /opt/pgsql
```

PostgreSQL root
postgres pghost1
postgres initdb

```
$ initdb -D /database/pg10/pg_root -E UTF8 --locale=C -U  
postgres -W
```

/database/pg10/pg_root
\$PGDATA/postgresql.conf

wal_level = replica	# minimal, replica, or
logical	
archive_mode = on	# enables archiving; off,
on, or always	
archive_command = '/bin/date'	# command to use to archive
a logfile segment	
max_wal_senders = 10	# max number of walsender
processes	
wal_keep_segments = 512	# in logfile segments, 16MB
each; 0 disables	
hot_standby = on	

postgres.conf

wal_level=WAL
minimal
replica
logical
minimal
WAL
replica
WAL
minimal
WAL
logical
WAL
10
WAL
replica
minimal
logical
WAL
replica
replica
replica

archive_mode=off
on
archive_command
WAL
on

archive_command=WAL
WAL
/bin/date
13

·max_wal_senders 参数控制WAL进程
pg_basebackup 参数控制WAL
max_connections 参数控制
10 参数控制WAL

·wal_keep_segments 参数控制pg_wal
WAL 参数控制WAL
WAL 参数控制WAL
WAL 参数控制WAL
16MB 参数控制--with-wal-segsize 参数控制WAL
pg_wal 参数控制
wal_keep_segments 参数控制x16MB 参数控制
512x16MB=8GB 参数控制pg_wal 参数控制WAL 参数控制

·hot_standby 参数控制
SQL 参数控制
on

postgresql.conf 参数控制
postgresql.conf 参数控制

pg_hba.conf 参数控制

```
# replication privilege.
host      replication      repuser      192.168.28.74/32
md5
```

```
host      replication    repuser    192.168.28.75/32
md5
```

```
pg_hba.conf
pg_hba.conf
```

```
pgghost1
```

```
[postgres@pgghost1 ~]$ pg_ctl start
```

```
postgres
repuser REPLICATION LOGIN
```

```
CREATE USER repuser
REPLICATION
LOGIN
CONNECTION LIMIT 5
ENCRYPTED PASSWORD 're12a345';
```

```
postgres
```

```
postgres=# SELECT pg_start_backup('francs_bk1');
pg_start_backup
```

```
0/4000060
(1 row)
```

```
pg_start_backup[REDACTED]
[REDACTED]pghost2[REDACTED]
```

```
$ tar czvf pg_root.tar.gz pg_root --exclude=pg_root/pg_wal
$ scp pg_root.tar.gz postgres@192.168.28.75:/database/pg10
```

[illegible]

pgghost2

```
$ tar xvf pg_root.tar.gz
```

```
postgres=# SELECT pg_stop_backup();
NOTICE:  pg_stop_backup complete, all required WAL segments
have been archived
      pg_stop_backup
-----
      0/2000130
(1 row)
```

pgghost2 recovery.conf
\$PGDATA
\$PGDATA

```
$ cp $PGHOME/share/recovery.conf.sample  
$PGDATA/recovery.conf
```

recovery.conf

```
recovery_target_timeline = 'latest'  
standby_mode = on  
primary_conninfo = 'host=192.168.28.74 port=1921  
user=repuser'
```

·recovery_target_timeline
timeline
latest
latest

·standby_mode
on
WAL
WAL

·primary_conninfo
IP
~/.pgpass

cd ~/.pgpass

```
[postgres@pghost2 ~]$ touch .pgpass
[postgres@pghost2 ~]$ chmod 0600 .pgpass
```

.pgpass 0600
pgpass

```
192.168.28.74:1921:replication:repuser:re12a345
192.168.28.75:1921:replication:repuser:re12a345
```

.pgpass IP
repuser
pghost2

```
$ pg_ctl start
```

WAL
WAL WAL
WAL

```
postgres 28575 28475  0 16:41 ?          00:00:00 postgres:
wal sender process repuser 192.168.28.75(57805) streaming
0/3025000
```

WAL

```
postgres 15449 15331 0 16:41 ? 00:00:00 postgres:
wal receiver process streaming 0/301FC68
```

[illegible]

```
postgres=# CREATE TABLE test_sr(id int4);
CREATE TABLE
postgres=# INSERT INTO test_sr VALUES (1);
INSERT 0 1
```

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □

```
postgres=# SELECT * FROM test_sr;
 id
----
  1
(1 row)
```

```

[postgres.conf] hot_standby [on]

```

```
hot_standby = on                # "off" disallows
queries during recovery
```

off postgresql

```
[postgres@pghost2 ~]$ psql postgres postgres
psql: FATAL:  the database system is starting up
```

```

pg_basebackup -h host -p 5432 -U postgres -D /var/lib/postgresql/12/base
$PGDATA/pg_log

```

12.1.2 pg_basebackup

```

pg_basebackup -h host -p 5432 -U postgres -D /var/lib/postgresql/12/base

```

```
1 pg_start_backup 'francs_bk1'
```

```
2 cp -r $PGDATA /var/lib/postgresql/12/base
```

```
3 pg_stop_backup
```

```

PostgreSQL
pg_basebackup -h host -p 5432 -U postgres -D /var/lib/postgresql/12/base
pg_start_backup pg_stop_backup
pg_basebackup

```

pg_basebackup 命令使用

```
pg_basebackup -D /database/pg10/pg_root
REPLICATION max_wal_senders
pg_basebackup -D /database/pg10/pg_root -WAL
pg_basebackup -D /database/pg10/pg_root
pgghost2
pg_basebackup -D /database/pg10/pg_root pgghost2
```

```
$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
$ rm -rf /database/pg10/pg_root
$ rm -rf /database/pg10/pg_tbs
```

pgghost2 pg_basebackup

```
$ pg_basebackup -D /database/pg10/pg_root -Fp -Xs -v -P -h
192.168.28.74 -p 1921
-U repuser
pg_basebackup: initiating base backup, waiting for
checkpoint to complete
pg_basebackup: checkpoint completed
pg_basebackup: write-ahead log start point: 1/B9000028 on
timeline 1
pg_basebackup: starting background WAL receiver
7791508/7791508 kB (100%), 2/2 tablespaces
pg_basebackup: write-ahead log end point: 1/B90039E0
pg_basebackup: waiting for background process to finish
```

```
streaming ...
pg_basebackup: base backup completed
```

```
pg_basebackup
checkpoint
$PGDATA
pg_basebackup
```

```
.-D
/database/pg10/pg_root
```

```
.-Fpg_basebackup
pplainttarpplain
$PGDATAttar
base.taroid.tarOID
OID
```

```
.-XWAL
ffetchsstreamffetch
WAL
wal_keep_segments
WAL
ffetchWALs
streamWAL
WALWAL
```

pg_basebackup --fetch --wal --wal-method=stream --no-wal-checkpoint --no-wal-archive --no-wal-rename --no-wal-copy --no-wal-delete --no-wal-delete-wal

pg_basebackup -v --verbose --no-wal-checkpoint --no-wal-archive --no-wal-rename --no-wal-copy --no-wal-delete --no-wal-delete-wal

pg_basebackup -P --progress --no-wal-checkpoint --no-wal-archive --no-wal-rename --no-wal-copy --no-wal-delete --no-wal-delete-wal

pg_basebackup -h --help -p --port -U --username --no-wal-checkpoint --no-wal-archive --no-wal-rename --no-wal-copy --no-wal-delete --no-wal-delete-wal
<https://www.postgresql.org/docs/10/static/app-pgbasebackup.html>

pg_basebackup --no-wal-checkpoint --no-wal-archive --no-wal-rename --no-wal-copy --no-wal-delete --no-wal-delete-wal
recovery.conf --no-wal-checkpoint --no-wal-archive --no-wal-rename --no-wal-copy --no-wal-delete --no-wal-delete-wal
\$PGDATA

```
$ cp ~/recovery.conf $PGDATA
```

pgghost2

```
$ pg_ctl start
```

pg_stat_replication WAL 12.1.3 pg_stat_replication sync_state

12.1.3 pg_stat_replication

pg_stat_replication pg_stat_replication sync_state

```
postgres=# SELECT
username,application_name,client_addr,sync_state
FROM pg_stat_replication ;
username | application_name | client_addr | sync_state
-----+-----+-----+-----
repuser  | walreceiver      | 192.168.28.75 | async
(1 row)
```

pg_stat_replication WAL sync_state

- async
- potential
- sync

·quorum PostgreSQL10
quorum standbys

sync_state async

12.2 配置

配置 PostgreSQL 的 WAL 模式。在 `postgresql.conf` 文件中，将 `wal_mode` 设置为 `replica`。这确保了 WAL 记录在写入主数据库时，也会同时写入到备库的 WAL 文件中。

配置 PostgreSQL 的 WAL 模式。在 `postgresql.conf` 文件中，将 `wal_mode` 设置为 `replica`。这确保了 WAL 记录在写入主数据库时，也会同时写入到备库的 WAL 文件中。

配置 PostgreSQL 的 WAL 模式。在 `postgresql.conf` 文件中，将 `wal_mode` 设置为 `replica`。这确保了 WAL 记录在写入主数据库时，也会同时写入到备库的 WAL 文件中。

12.2.1 synchronous_commit 配置

配置 PostgreSQL 的 `synchronous_commit` 参数。在 `postgresql.conf` 文件中，将 `synchronous_commit` 设置为 `on`。这确保了在提交事务时，所有副本的 WAL 记录都已写入，从而保证了数据的一致性。

配置 PostgreSQL 的 `synchronous_commit` 参数。在 `postgresql.conf` 文件中，将 `synchronous_commit` 设置为 `on`。这确保了在提交事务时，所有副本的 WAL 记录都已写入，从而保证了数据的一致性。

WAL
on off local remote_apply
remote_write

·on WAL WAL BUFFER
WAL on WAL
WAL on

·off off WAL BUFFER
WAL off
off
off
off

·local local on
WAL

·remote_write
WAL
WAL
WAL WAL

```

-- WAL remote_write
WAL WAL
WAL

```

```

    ·on on
WAL WAL
on WAL WAL
WAL

```

```

    ·remote_apply
WAL WAL
remote_apply WAL
WAL
WAL

```

12.2.2

recovery.conf

```
primary_conninfo = 'host=192.168.28.74 port=1921
user=repuser application_name=node2'
```

```

primary_conninfo
application_name application_name
postgresql.conf

```

```
synchronous_standby_names=''
application_name=''node2'
```

```
postgresql.conf

```

```
synchronous_commit = on
synchronous_standby_names = 'node2'
```

```
wal_level=logical

```

```
·synchronous_commit=on
remote_apply=on

```

```
·synchronous_standby_names=
node2
recovery.conf
primary_conninfo=application_name=

```

```


```

```
[postgres@pghost1 ~]$ pg_ctl reload
server signaled
```

```
wal_level[]
synchronous_commit[]
synchronous_standby_names[]
[]pg_ctl reload[]
[]wal_level[]replica[]
[]
```

```
[]recovery.conf[]pghost2
[]
```

```
[postgres@pghost2 ~]$ pg_ctl restart -m fast
waiting for server to shut down.... done
server stopped
```

```
[]WAL[]WAL[]
[]
```

```
[]
```

```
postgres=# SELECT
username,application_name,client_addr,sync_state
FROM pg_stat_replication ;
username | application_name | client_addr | sync_state
-----+-----+-----+-----
repuser  | node2            | 192.168.28.75 | sync
(1 row)
```

```

pg_stat_replication sync_state
sync sync

```

12.2.3 字符串“ ”

[illegible]

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □

```
[postgres@pghost2 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □

```
postgres=# SELECT * FROM test_sr LIMIT 1;
 id
----
  1
(1 row)
```

データベースのインストールと起動

データベースのインストールと起動

```
postgres=# INSERT INTO test_sr(id) VALUES (5);  
--
```

データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動

データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動
データベースのインストールと起動
12.7 データベースのインストールと起動

12.3 数据库性能测试

PostgreSQL数据库性能测试工具pgbench，可以在一台机器上模拟多用户并发访问数据库，进行各种类型的数据库操作（比如连接、插入、更新、查询等），以用来测试数据库在不同负载情况下的性能。pgbench默认使用4个CPU，8GB内存。

修改postgresql.conf配置文件

```
wal_level = replica          # minimal, replica, or
logical                     # synchronization level;
synchronous_commit = off
```

修改postgresql.conf配置文件

```
wal_level = replica          # minimal, replica, or
logical                     # synchronization level;
synchronous_commit = off    # in logfile segments, 16MB
wal_keep_segments = 512    # each; 0 disables
```

修改postgresql.conf配置文件

```
wal_level = replica          # minimal, replica, or logical
synchronous_commit = on      # synchronization level;
synchronous_standby_names = 'node2' # standby servers
that provide sync rep
```

```
wal_keep_segments = 512      # in logfile segments, 16MB
each; 0 disables
```

```

postgresql.conf postgresql.conf

```

12.3.1

```

test_per1 1000

```

```
postgres=# CREATE TABLE test_per1(
    id int4,
    name text,
    create_time timestamp(0) without time zone default
clock_timestamp());
CREATE TABLE
postgres=# INSERT INTO test_per1(id,name)
    SELECT n,n||'_per1'
    FROM generate_series(1,10000000) n;
INSERT 0 10000000
```

```
postgres=# ALTER TABLE test_per1 ADD PRIMARY KEY(id);
ALTER TABLE
postgres=# ANALYZE test_per1;
ANALYZE
```

SQL select_per1.sql
1000000

```
\set v_id random(1,100000000)
```

```
SELECT name FROM test_per1 WHERE id=:v_id;
```

v_id 1 1000
test_per1 2 4 8 16 TPS
pgbench

```
pgbench -c 2 -T 120 -d postgres -U postgres -n N -M prepared  
-f select_per1.sql > select_2.out 2>&1 &  
pgbench -c 4 -T 120 -d postgres -U postgres -n N -M prepared  
-f select_per1.sql > select_4.out 2>&1 &  
pgbench -c 8 -T 120 -d postgres -U postgres -n N -M prepared  
-f select_per1.sql > select_8.out 2>&1 &  
pgbench -c 16 -T 120 -d postgres -U postgres -n N -M  
prepared -f select_per1.sql > select_16.out 2>&1 &
```

pgbench 120 -M prepared
prepared statements -n VACUUM
pgbench
12-2

12-2

连接数	单实例 (TPS)	异步流复制 (TPS)	同步流复制 (TPS)
2	17061	16659	16427
4	20370	19966	19622
8	18193	16756	17991
16	11707	11745	11675

图 12-1 不同连接数下的 TPS 性能对比

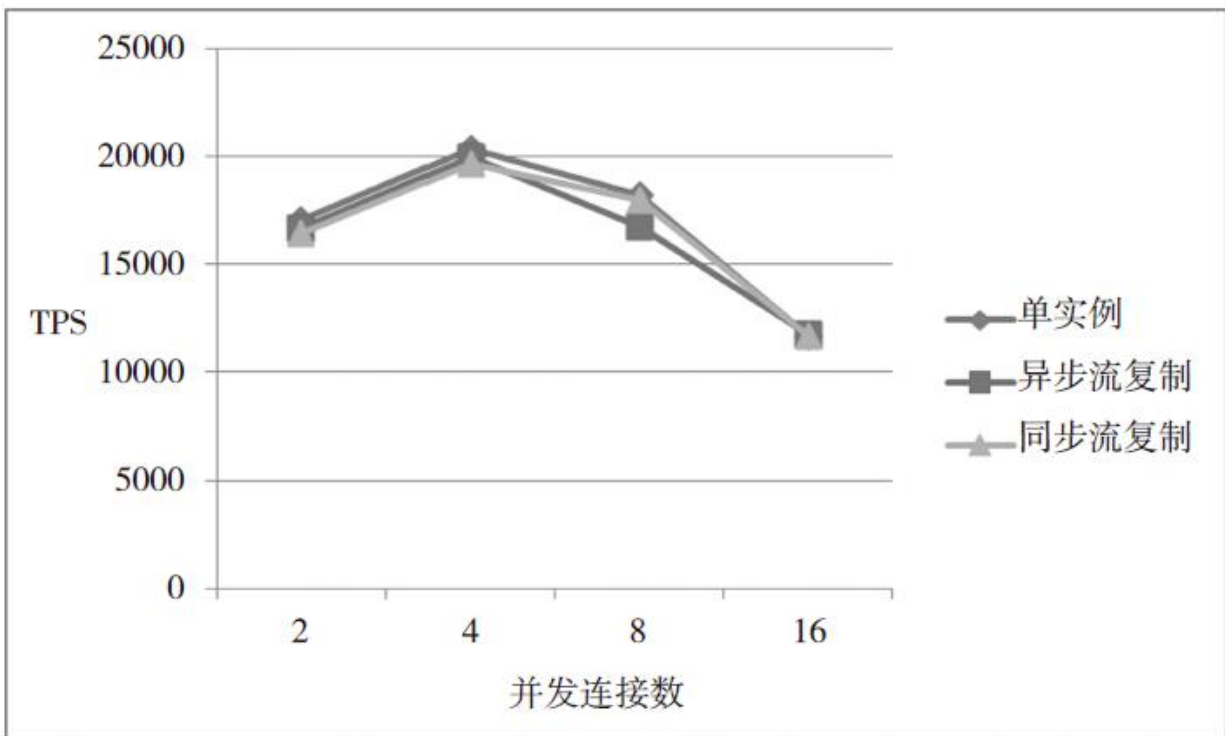


图 12-1 不同连接数下的 TPS 性能对比

图 12-1 展示了在不同连接数下，单实例、异步流复制和同步流复制三种配置的性能对比。从图中可以看出，随着连接数的增加，TPS 性能会先上升后下降。在 4 个连接数时，性能达到峰值。在 16 个连接数时，性能显著下降。此外，同步流复制的性能在 4 个连接数时略高于异步流复制，但在 16 个连接数时，两者的性能差距缩小。

图 12-1 展示了在不同连接数下，单实例、异步流复制和同步流复制三种配置的性能对比。从图中可以看出，随着连接数的增加，TPS 性能会先上升后下降。在 4 个连接数时，性能达到峰值。在 16 个连接数时，性能显著下降。此外，同步流复制的性能在 4 个连接数时略高于异步流复制，但在 16 个连接数时，两者的性能差距缩小。



CPU使用率を50%程度に抑え、pgbenchのテスト結果を比較し、tpsを

12.3.2 テスト環境

テスト環境として、test_per2というテーブルを作成し、1000万行のデータを
 挿入する。

```

postgres=# CREATE TABLE test_per2(id int4,name text,flag
char(1));
CREATE TABLE
postgres=# INSERT INTO test_per2(id,name)
SELECT n,n||'_per2'
FROM generate_series(1,10000000) n;
INSERT 0 10000000
  
```

テスト環境として、test_per2というテーブルを作成し、1000万行のデータを

```

postgres=# ALTER TABLE test_per2 ADD PRIMARY KEY(id);
ALTER TABLE
postgres=# ANALYZE test_per2;
ANALYZE
  
```

テスト環境として、test_per2というテーブルを作成し、1000万行のデータを
 挿入する。

```
\set v id random(1,1000000)
```

```

v_id11000
test_per2flag24816
TPSpgbench

```

pgbench120pgbench
12-3

连 接 数	单实例 (TPS)	异步流复制 (TPS)	同步流复制 (TPS)
2	13062	13640	856
4	17306	17258	1270
8	10715	10598	2414
16	8793	9388	4681

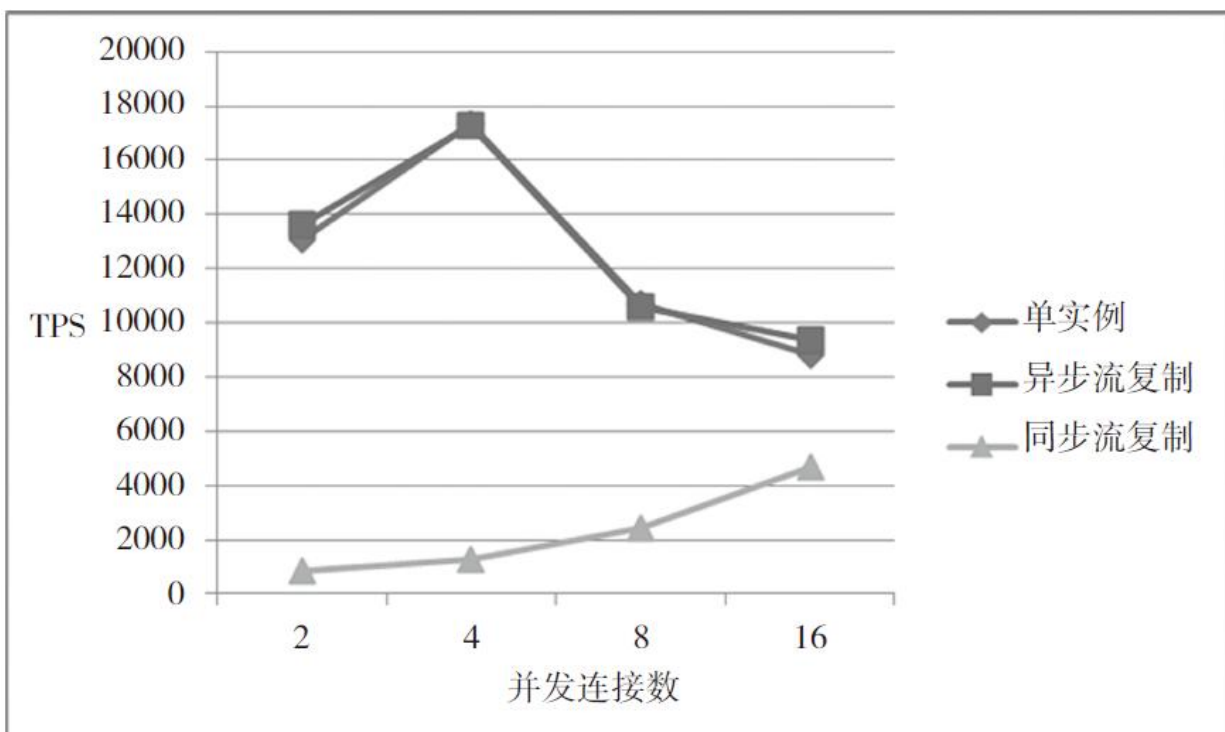


图12-2 不同并行连接数下的性能对比

该图展示了在不同并行连接数下，三种复制方式的性能对比。Y轴为TPS（每秒事务数），X轴为并行连接数。异步流复制（异步流复制）在4个并行连接时达到最高性能，约为17500 TPS。同步流复制（同步流复制）随着并行连接数的增加，性能持续提升，在16个并行连接时达到约4800 TPS。单实例（单实例）在2个并行连接时性能约为13000 TPS，但在16个并行连接时性能下降至约9000 TPS。

12.4 复制

复制是数据库系统的重要功能之一，它允许用户将数据从一个数据库复制到另一个数据库，或者将数据从一个数据库复制到另一个数据库的副本中。

12.4.1 pg_stat_replication

pg_stat_replication 视图显示了当前正在运行的复制进程的状态。该视图包含以下列：

```
postgres=# SELECT * FROM pg_stat_replication ;
-[ RECORD 1 ]-----+-----
pid                | 7683
usesysid           | 16384
username           | repuser
application_name    | node2
client_addr        | 192.168.28.75
client_hostname     |
client_port        | 57870
backend_start       | 2017-09-05 11:50:31.629468+08
backend_xmin        |
state               | streaming
sent_lsn            | 3/643CB568
write_lsn           | 3/643CB568
flush_lsn           | 3/643CB488
replay_lsn          | 3/643CB030
write_lag           | 00:00:00.000224
flush_lag           | 00:00:00.001562
replay_lag          | 00:00:00.006596
sync_priority       | 1
sync_state          | sync
```

□□□□□□□□□□□□

·pid□WAL□□□□□□□□

·username□WAL□□□□□□□□□□

·application_name□□WAL□□□□□□□□□□
□□□□□□□□recovery.conf□□□□□
primary_conninfo□□application_name□□□
□□

·client_addr□□□□WAL□□□□□□□□IP□□□□□
□□□□IP□

·backend_start□WAL□□□□□□□□□□

·state□□□WAL□□□□□□□□startup□□WAL□
□□□□□□□□catchup□□□□□□□□□□streaming□
□□□□□□□□□□□□□□□□□□□□WAL□□□□□□□□□□
□□□□□□□□backup□□□□pg_basebackup□□□□□□
□□stopping□□WAL□□□□□□□□□□

·sent_lsn□WAL□□□□□□□□□□WAL□□□□□

·write_lsn□□□□□□□□□□WAL□□□□□□□□WAL□□
□□□□□□□□□□□□□□□□□□□□WAL□□□□□

·flush_lsn WAL WAL
WAL

·replay_lsn WAL

·write_lag WAL WAL
WAL WAL WAL
WAL

·flush_lag WAL WAL
WAL WAL WAL WAL
WAL

·replay_lag WAL WAL
WAL WAL WAL WAL
WAL

·sync_priority
quorum

·sync_state async
potential
sync
quorum quorum standbys
quorum standbys

write_lag flush_lag replay_lag
PostgreSQL10

12.4.2

WAL
WAL

WAL

WAL write flush replay
pg_stat_replication write_lag
flush_lag replay_lag
WAL
SQL

```
postgres=# SELECT
pid,username,client_addr,state,write_lag,flush_lag,replay_lag
FROM pg_stat_replication ;
-[ RECORD 1 ]-----
pid          | 7683
username     | repuser
client_addr  | 192.168.28.75
state        | streaming
write_lag    | 00:00:00.000997
flush_lag    | 00:00:00.002008
replay_lag   | 00:00:00.002916
```

WAL
WAL
WAL write_lag
flush_lag replay_lag

replay_lag > flush_lag > write_lag

flush_lag 0.2008
replay_lag 0.2916 replay_lag
flush_lag WAL WAL
WAL replay_lag flush_lag

write_lag flush_lag replay_lag
PostgreSQL 10 10
pg_stat_replication
SQL

```
postgres=# SELECT EXTRACT(SECOND FROM now() -
pg_last_xact_replay_timestamp());
      date_part
-----
0.002227
(1 row)
```

pg_last_xact_replay_timestamp
WAL

```

pg_last_xact_replay_timestamp()
--
          WAL replay timestamp
          pg_last_xact_replay_timestamp()
          returns the timestamp of the last xact
          replayed from the WAL.

```

WAL

WAL WAL WAL SQL

```
postgres=# SELECT pid,username,client_addr,state,
                pg_wal_lsn_diff(pg_current_wal_lsn(),write_lsn)
write_delay,
                pg_wal_lsn_diff(pg_current_wal_lsn(),flush_lsn)
flush_delay,
                pg_wal_lsn_diff(pg_current_wal_lsn(),replay_lsn)
replay_dely
        FROM pg_stat_replication ;
-[ RECORD 1 ]-----
pid          | 7683
username     | repuser
client_addr  | 192.168.28.75
state        | streaming
write_delay  | 560
flush_delay  | 896
replay_dely  | 1272
```

```
pg_current_wal_lsn WAL
pg_wal_lsn_diff WAL
WAL write
560 flush 896 replay 1272
```

pg_stat_wal_receiver

pg_stat_wal_receiver
pg_stat_wal_receiver
pg_stat_wal_receiver
pg_stat_wal_receiver

12.4.3 pg_stat_wal_receiver

pg_stat_replication WAL
WAL

```
postgres=# SELECT * FROM pg_stat_wal_receiver ;
-[ RECORD 1 ]-----+-----
pid                | 22573
status             | streaming
receive_start_lsn  | 3/2D000000
receive_start_tli  | 1
received_lsn       | 3/852DC428
received_tli       | 1
last_msg_send_time | 2017-09-06 15:35:28.178167+08
last_msg_receipt_time | 2017-09-06 15:35:28.177706+08
latest_end_lsn     | 3/852DC508
latest_end_time    | 2017-09-06 15:35:28.178167+08
slot_name          |
conninfo           | user=repuser
passfile=/home/postgres/.pgpass dbname=replication
host=192.168.28.74 port=1921 application_name=node2
fallback_application_name=walreceiver sslmode=disable
sslcompression=1 target_session_attrs=any
```

pg_stat_wal_receiver

- pid WAL
- status WAL
- receive_start_lsn WAL
WAL
- received_lsn WAL WAL
- last_msg_send_time
- last_msg_receipt_time
- conninfo WAL
\$PGDATA/recovery.conf
primary_conninfo

12.4.4

PostgreSQL WAL

2017-10-07 09:04:59.67741+08
(1 row)

WAL

```
postgres=# SELECT pg_current_wal_lsn();
           pg_current_wal_lsn
-----
          3/940001B0
(1 row)
```

WAL

```
postgres=# SELECT
pg_wal_lsn_diff('3/940001B0','3/940001A0');
           pg_wal_lsn_diff
-----
                  16
(1 row)
```

12.5 主备切换

主备切换是指将主库的日志流复制到备库，当主库发生故障时，备库可以接管主库的工作。在 PostgreSQL 中，主备切换可以通过 `pg_ctl promote` 命令来实现。该命令用于将备库提升为主库，并使其开始接收新的日志流。

在 PostgreSQL 中，主备切换可以通过 `pg_ctl promote` 命令来实现。该命令用于将备库提升为主库，并使其开始接收新的日志流。

12.5.1 主备切换的步骤

主备切换的步骤如下：

1. 确保主库和备库的配置一致。
2. 确保主库和备库的日志流处于同步状态。
3. 在主库上执行 `pg_ctl promote` 命令，将主库提升为主库。
4. 在备库上执行 `pg_ctl promote` 命令，将备库提升为主库。

主备切换的步骤如下：

主备切换的步骤如下：

1. 在主库上执行 `pg_ctl promote` 命令，将主库提升为主库。
2. 在备库上执行 `pg_ctl promote` 命令，将备库提升为主库。
3. 在主库上执行 `wal sender..streaming` 命令，将主库的日志流复制到备库。

```
[postgres@pghost1 ~]$ ps -ef | grep "wal" | grep -v "grep"
postgres 16666 16661  0 Sep06 ?                00:00:09 postgres: wal
writer process
postgres 16672 16661  0 Sep06 ?                00:00:13 postgres: wal
sender process repuser 192.168.28.75(57872) streaming
3/9C34BCB8
```

wal receiver..streaming

```
[postgres@pghost2 ~]$ ps -ef | grep "wal" | grep -v "grep"
postgres 27291 22567  0 Sep06 ?                00:00:32 postgres: wal
receiver process      streaming 3/9C355788
```

WAL WAL

WAL WAL
pg_stat_replication

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state
FROM pg_stat_replication ;
 pid | username | application_name | client_addr |
state | sync_state
-----+-----+-----+-----+-----
 16672 | repuser  | node2            | 192.168.28.75 |
streaming | sync
(1 row)
```

pg_stat_wal_receiver

```
postgres=# SELECT
pid,status,last_msg_send_time,last_msg_receipt_time,conninfo
FROM pg_stat_wal_receiver ;
```

```
-[ RECORD1 ]-----+-----
pid          | 17551
status       | streaming
last_msg_send_time | 2017-10-07 09:22:07.479282+08
last_msg_receipt_time | 2017-10-07 09:22:07.480277+08
conninfo     | user=repuser
passfile=/home/postgres/.pgpass dbname= replication
host=192.168.28.74 port=1921 application_name=slavel
fallback_application_name=walreceiver sslmode=disable
sslcompression=1 target_session_attrs=any
```

```
postgres=# SELECT pg_is_in_recovery();
pg_is_in_recovery
```

```
-----
t
(1 row)
```

pg_controldata WAL checkpoint Database cluster state

```
[postgres@pghost1 ~]$ pg_controldata | grep cluster
Database cluster state:          in production
```

in production in archive recovery

```
[postgres@pghost2 ~]$ pg_controldata | grep cluster
Database cluster state:          in archive recovery
```

recovery.conf

\$PGDATA recovery.conf \$PGDATA recovery.done

12.5.2

PostgreSQL9.0 pghost1 pghost2

1 recovery.conf trigger_file
recovery.conf

2 -m fast

3 recovery.conf
recovery.done

4 \$PGDATA
recovery.conf
recovery.conf
\$PGHOME/recovery.conf.sample
recovery.done recovery.done
recovery.conf
primary_conninfo IP IP

5
recovery.conf

recovery.conf

```
recovery_target_timeline = 'latest'
standby_mode = on
primary_conninfo = 'host=192.168.28.74 port=1921
user=repuser'
trigger_file =
'/database/pg10/pg_root/.postgresql.trigger.1921'
```

trigger_file
[...]

[...]

```
[postgres@pghost1 ~]$ pg_ctl stop -m fast
waiting for server to shut down... done
server stopped
```

[...]

```
[postgres@pghost2 pg_root]$ ll recovery.conf
-rw-r--r-- 1 postgres postgres 5.8K Sep  8 20:47
recovery.conf
[postgres@pghost2 pg_root]$ touch
/database/pg10/pg_root/.postgresql.trigger.1921
```

recovery.conf
trigger_file
recovery
.conf.done

```
[postgres@pghost2 pg_root]$ ll recovery.done
-rw-r--r-- 1 postgres postgres 5.8K Sep  8 20:47
recovery.done
```

[...]

2017-09-08 21:00:21.622 CST,,,4357,,59b29465.1105,1,,2017-09-08 21:00:21 CST,,0,FATAL,XX000,"could not connect to the

primary server: could not connect to server: Connection refused

Is the server running on host "192.168.28.74" and accepting

TCP/IP connections on port 1921?,,,,,,,,,""
2017-09-08 21:00:26.622 CST,,,4235,,59b2916f.108b,9,,2017-09-08 20:47:43 CST,1/0,0,LOG,00000,"trigger file found: /database/pg10/pg_root/.postgresql.trigger.1921",,,,,,,,,,""
2017-09-08 21:00:26.622 CST,,,4235,,59b2916f.108b,10,,2017-09-08 20:47:43 CST,1/0,0,LOG,00000,"redo done at 3/A0000028",,,,,,,,,,""
2017-09-08 21:00:26.622 CST,,,4235,,59b2916f.108b,11,,2017-09-08 20:47:43 CST,1/0,0,LOG,00000,"last completed transaction was at log time 2017-09-08 20:59:30.746045+08",,,,,,,,,,""
2017-09-08 21:00:26.640 CST,,,4235,,59b2916f.108b,12,,2017-09-08 20:47:43 CST,1/0,0,LOG,00000,"selected new timeline ID: 2",,,,,,,,,,""
2017-09-08 21:00:26.792 CST,,,4235,,59b2916f.108b,13,,2017-09-08 20:47:43 CST,1/0,0,LOG,00000,"archive recovery complete",,,,,,,,,,""
2017-09-08 21:00:26.808 CST,,,4233,,59b2916f.1089,3,,2017-09-08 20:47:43 CST,,0,LOG,00000,"database system is ready to accept connections",,,,,,,,,,""

XX
XX

XXXXpg_controldataXXXXXXXXXXXXXXXXXXXX

[postgres@pghost2 pg_root]\$ pg_controldata | grep cluster
Database cluster state: in production

XX
test_alivedXXXXXXXXXXXXXXXXXXXX

```
postgres=# CREATE TABLE test_alived2(id int4);
CREATE TABLE
postgres=# INSERT INTO test_alived2 VALUES(1);
INSERT 0 1
```

4 recovery.conf

```
recovery_target_timeline = 'latest'
standby_mode = on
primary_conninfo = 'host=192.168.28.75 port=1921
user=repuser'
trigger_file =
'/database/pg10/pg_root/.postgresql.trigger.1921'
```

pghost2 recovery.done
primary_conninfo host IP

pghost1 postgres
~/.pgpass

```
[postgres@pghost1 ~]$ touch ~/.pgpass
[postgres@pghost1 ~]$ chmod 600 ~/.pgpass
```

~/.pgpass

```
192.168.28.74:1921:replication:repuser:re12a345
192.168.28.75:1921:replication:repuser:re12a345
```

```
[postgres@pghost1 ~]$ pg_ctl start
```

```
postgres=# SELECT * from test_alived2 ;
   id
-----
   1
(1 row)
```

[illegible]

12.5.3 pg_ctl promote

PostgreSQL 9.1
pg_ctl promote
promote

```
pg_ctl promote [-D datadir]
```

-D \$PGDATA
promote

pg_ctl promote
1 recovery.conf trigger_file
3 pg_ctl promote

1 -m fast

2 pg_ctl promote
recovery.conf recovery.done

3 \$PGDATA
recovery.conf
recovery.conf
\$PGHOME/recovery.conf.sample

pg_ctl promote
waiting for server to promote.... done
server promoted

pg_ctl promote
waiting for server to promote.... done
server promoted

pg_ctl promote
waiting for server to promote.... done
server promoted

pg_ctl promote
waiting for server to promote.... done
server promoted

pg_ctl promote
waiting for server to promote.... done
server promoted

pg_ctl promote
waiting for server to promote.... done
server promoted

```
[postgres@pghost1 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

```
$PGDATA/recovery.done
recovery.conf
```

```
[postgres@pghost1 pg_root]$ mv recovery.done recovery.conf
```

```
pghost1/recovery.conf
```

```
recovery_target_timeline = 'latest'
standby_mode = on
primary_conninfo = 'host=192.168.28.75 port=1921
user=repuser'
```

```
pghost1
```

```
[postgres@pghost1 pg_root]$ pg_ctl start
```

```
pghost1/WAL
pghost2/WALpghost1

```

```
2017-09-09 13:22:45.540 CST,,,11948,,59b37aa5.2eac,2,,2017-
09-09 13:22:45 CST,,0,FATAL,XX000,"could not start WAL
streaming: ERROR: requested starting point 3/A3000000 on
timeline 3 is not in this server's history
```

DETAIL: This server's history forked from timeline 3 at
3/A213F930.",,,,,,,,,""
2017-09-09 13:22:45.540 CST,,,11944,,59b37aa5.2ea8,5,,2017-
09-09 13:22:45 CST,1/0,0,LOG,00000,"new timeline 4 forked off
current database system timeline 3 before current recovery
point 3/A3000098",,,,,,,,,""
2017-09-09 13:22:45.543 CST,,,11949,,59b37aa5.2ead,1,,2017-
09-09 13:22:45 CST,,0,FATAL,XX000,"could not start WAL
streaming: ERROR: requested starting point 3/A3000000 on
timeline 3 is not in this server's history
DETAIL: This server's history forked from timeline 3 at
3/A213F930.",,,,,,,,,""
2017-09-09 13:22:45.543 CST,,,11944,,59b37aa5.2ea8,6,,2017-
09-09 13:22:45 CST,1/0,0,LOG,00000,"new timeline 4 forked off
current database system timeline 3 before current recovery
point 3/A3000098",,,,,,,,,""

pgghost1
pgghost1

PostgreSQL pg_rewind
pg_basebackup pg_rewind
pgghost1 pgghost2

pg_rewind

postgresql.conf wal_log_hints
on

·`initdb --data-checksums` `I/O`
`initdb`

`initdb --data-checksums`
`postgresql.conf` `wal_log_hints`

```
wal_log_hints = on
```

`pghost2`

```
[postgres@pghost2 pg_root]$ pg_ctl promote
waiting for server to promote.... done
server promoted
```

`pghost1`
`pghost1`
`pghost1`

```
[postgres@pghost1 pg_root]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

pg_rewindpgghost2
pgghost1

```
[postgres@pgghost1 pg_root]$ pg_rewind --target-pgdata $PGDATA  
--source-server= 'host= 192.168.28.75 port=1921 user=postgres  
dbname=postgres' -P  
connected to server  
servers diverged at WAL location 3/A7006508 on timeline 5  
rewinding from last common checkpoint at 3/A7000028 on  
timeline 5  
reading source file list  
reading target file list  
reading WAL in target  
need to copy 7663 MB (total source directory size is 9237 MB)  
7847309/7847309 kB (100%) copied  
creating backup label and updating control file  
syncing target data directory  
Done!
```

postgres
postgres ~/.pgpass

recovery.done recovery.conf

```
[postgres@pgghost1 pg_root]$ mv recovery.done recovery.conf
```

recovery.conf primary_conninfo
host

```
[postgres@pgghost1 pg_root]$ pg_ctl start
```

pgghost1
pgghost1

12.6 备份

数据库备份是数据库管理中最重要的一环。本章将介绍如何使用pg_dump和pg_dumpall工具进行数据库备份。pg_dumpall工具用于备份整个数据库集群，而pg_dump工具用于备份指定的数据库。本章将详细介绍pg_dumpall工具的使用，包括如何备份整个数据库集群以及如何备份指定的数据库。

12.6.1 备份整个数据库集群

PostgreSQL数据库集群的备份可以通过pg_dumpall工具实现。pg_dumpall工具位于/usr/bin目录下。使用pg_dumpall工具备份整个数据库集群的命令如下：

```
pg_dumpall -h host -U user > backup.sql
```

其中，host表示数据库集群的主机名，user表示数据库用户。备份文件backup.sql将存储在当前目录下。pg_dumpall工具支持多种输出格式，包括文本格式、二进制格式和目录格式。本章将详细介绍pg_dumpall工具的使用，包括如何备份整个数据库集群以及如何备份指定的数据库。

12.6.2 备份指定的数据库

pg_dumpall工具可以用于备份整个数据库集群，而pg_dump工具可以用于备份指定的数据库。pg_dump工具位于/usr/bin目录下。使用pg_dump工具备份指定的数据库的命令如下：

```
pg_dump -h host -U user database > backup.sql
```

其中，host表示数据库集群的主机名，user表示数据库用户，database表示要备份的数据库名。备份文件backup.sql将存储在当前目录下。pg_dump工具支持多种输出格式，包括文本格式、二进制格式和目录格式。本章将详细介绍pg_dump工具的使用，包括如何备份指定的数据库以及如何备份指定的表。

recovery_min_apply_delay
recovery.conf

recovery_min_apply_delay (integer)

·ms

·s

·min

·h

·d

WAL
WAL
WAL
WAL
WAL
WAL

pghost21

```
recovery_min_apply_delay = 1min
```

```

      1

```

```
[postgres@pghost2 pg_root]$ pg_ctl restart
```

```
test_delay
```

```
postgres=# CREATE TABLE test_delay(id int4,create_time
timestamp(0) without time zone);
CREATE TABLE
```

```


```

```
postgres=# INSERT INTO test_delay(id,create_time) VALUES
(1,now());
INSERT 0 1
```

```
test_delay
SQL
```

```
postgres=# SELECT now(),create_time FROM test_delay;
-[ RECORD 1 ]-----
now          | 2017-09-10 16:18:50.414074+08
create_time  | 2017-09-10 16:17:50
```

pg_dump -c -Ft -f test_delay.sql test_delay

pg_restore -c -Ft -f test_delay.sql test_delay


pg_dumpall -c -Ft -f test_delay.sql

```
postgres=# DROP TABLE test_delay ;
DROP TABLE
```

pg_dumpall -c -Ft -f test_delay.sql

```
postgres=# SELECT * FROM test_delay;
 id |      create_time
-----+-----
   1 | 2017-09-10 16:17:50
(1 row)
```

pg_dumpall -c -Ft -f test_delay.sql

 recovery_min_apply_delay
pg_wal WAL

pg_wal
pg_wal

12.6.3 recovery_min_apply_delay

recovery_min_apply_delay
synchronous_commit on
remote_apply on
WAL WAL
remote_apply
WAL WAL WAL
12.2.1

synchronous_commit
on remote_apply

synchronous_commit on
recovery_min_apply_delay 1

synchronous_commit
pg_ctlreload
recovery_min_apply_delay
pg_ctlreload

test_delay
test_delay

```
postgres=# INSERT INTO test_delay(id,create_time)
VALUES(1,now());
INSERT 0 1
```

test_delay

```
postgres=# SELECT now(),create_time FROM test_delay;
-[ RECORD 1 ]-----
now          | 2017-09-10 16:58:22.526087+08
create_time  | 2017-09-10 16:57:22
```

synchronous_commit
on

synchronous_commit
remote_apply
recovery_min_apply_delay1

SQLtest_delay

```
postgres=# INSERT INTO test_delay(id,create_time)
VALUES(2,now());
--
```

SQL SQL

```
postgres=# \timing
Timing is on.
postgres=# INSERT INTO test_delay(id,create_time)
VALUES(3,now());
INSERT 0 1
Time: 60008.295 ms (01:00.008)
```

SQL 60 INSERT

60 synchronous_commit

remote_apply INSERT

INSERT WAL WAL

WAL

WAL

synchronous_commit on

synchronous_commit remote_apply

recovery_min_apply_delay

synchronous_commit remote_apply

12.7 配置高可用

配置高可用需要配置主备节点，主节点配置如下：
PostgreSQL 配置主节点，主节点配置如下：
配置 PostgreSQL 9.6 主节点，主节点配置如下：
PostgreSQL 10 主节点
synchronous_standby_names=ANY
配置主节点，主节点配置如下：
Quorum
PostgreSQL 10 主节点配置如下：
配置主节点，主节点配置如下：

配置主节点，主节点配置如下：
pghost3
12-4

12-4 配置高可用

主 机	主 机 名	IP 地址	操 作 系 统	PostgreSQL 版本
主节点	pghost1	192.168.28.74	CentOS6.9	PostgreSQL10
备节点 1	pghost2	192.168.28.75	CentOS6.9	PostgreSQL10
备节点 2	pghost3	192.168.28.76	CentOS6.9	PostgreSQL10

配置高可用需要配置主备节点，主节点配置如下：
pg_stat_replication

```
postgres=# SELECT  
pid,username,application_name,client_addr,state,sync_state,sy  
nc_priority  
FROM pg_stat_replication ;
```

pid	username	application_name	client_addr
state	sync_state	sync_priority	
26030	repuser	node2	192.168.28.75
streaming	async		0
2799	repuser	node3	192.168.28.76
streaming	async		0

(2 rows)

node2 node3
sync_state

synchronous_standby_names

12.7.1 synchronous_standby_names

synchronous_standby_names string
PostgreSQL 10

·standby_name[...]

·[FIRST]num_sync
standby_name[...]

·ANY num_sync standby_name[...]

同步
synchronous_standby_names=standby_name[...]

standby_name 指定 standby 数据库名称
\$PGDATA/recovery.conf 文件
primary_conninfo 指定 application_name 名称
指定 standby 数据库名称
指定 standby 数据库名称 9.5 版本 9.5 版本

指定 standby 数据库名称 's1 s2' 指定 standby 数据库名称 s1 s2
指定 standby 数据库名称 s1 s2

同步 synchronous_standby_names=
[FIRST] num_sync standby_name[...]

FIRST 指定 standby 数据库名称
指定 standby 数据库名称 num_sync
指定 standby 数据库名称

```
synchronous_standby_names = 'FIRST 2(s1,s2,s3)'
```

指定 standby 数据库名称 s1 s2 指定 standby 数据库名称 s1 s2
指定 standby 数据库名称 s1 s2
指定 standby 数据库名称 WAL 指定 standby 数据库名称 s3
指定 standby 数据库名称 s1 s2 s3

同步
synchronous_standby_names=ANY
num_sync[standby_name[...]]

ANY 同步 quorum 同步
num_sync 同步 s1 s2 s3 s4
同步

```
synchronous_standby_names = 'ANY 2 (s1, s2, s3)'
```

ANY 2 同步
同步 WAL 同步 WAL
同步 s4 同步 s4

同步 Quorum 同步

12.7.2 同步

pghost1 postgresql.conf

```
synchronous_standby_names = 'first 1 (node2,node3)'
```

first 1 同步
同步 node2 node3
同步

pg_ctl reload

```
[postgres@pghost1 pg_root]$ pg_ctl reload
server signaled
```

pg_ctl reload

```
postgres=# show synchronous_standby_names ;
synchronous_standby_names
-----
first 1 (node2,node3)
(1 row)
```

pg_stat_replication

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sync_priority FROM pg_stat_replication ;
 pid | username | application_name | client_addr | state | sync_state | sync_priority
-----+-----+-----+-----+-----+-----+-----
 1536 | repuser  | node2            | 192.168.28.75 | streaming | sync | 1
 2799 | repuser  | node3            | 192.168.28.76 | streaming | potential | 2
(2 rows)
```

sync_state node2
sync_state async sync node2
sync_priority 1 node3 sync_state

```
async potential node3
sync_priority 2 sync
potential
```

node2

```
[postgres@pghost2 pg_root]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

pg_stat_replication

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sync_priority FROM pg_stat_replication ;
   pid   | username | application_name | client_addr | state | sync_state | sync_priority |
-----+-----+-----+-----+-----+-----+-----
    2799 | repuser  | node3            | 192.168.28.76 | streaming | sync      | 2              |
(1 row)
```

node3 sync_state sync

```
postgres=# INSERT INTO test_delay(id,create_time) VALUES
(4,now());
INSERT 0 1
```

node2 和 node3 的配置文件

```
[postgres@pghost3 pg_root]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

配置主节点

```
postgres=# INSERT INTO test_delay(id,create_time) VALUES
(5,now());
-- 插入成功
```

配置从节点

配置主节点

```
synchronous_standby_names='first
1 node2 node3'
```

配置从节点

```
node2:
primary_conninfo='host=192.168.1.1 port=5432 user=postgres password=postgres'
node3:
primary_conninfo='host=192.168.1.1 port=5432 user=postgres password=postgres'
```

12.7.3 Quorum

Quorum 是 PostgreSQL 10 引入的一个新特性，它允许在同步复制模式下，主节点在提交事务时，只需要等待任意一个从节点响应即可。


```

-----+-----+-----+-----+-----+-----+
-----+-----
      26906 | repuser | node3          | 192.168.28.76 |
streaming| quorum   |                | 1              |
      26926 | repuser | node2          | 192.168.28.75 |
streaming| quorum   |                | 1              |
(2 rows)

```

```

node2 node3 sync_state
quorum sync_priority 1
Quorum sync_priority
node2

```

```

[postgres@pghost2 pg_root]$ pg_ctl stop
waiting for server to shut down.... done
server stopped

```

```

pg_stat_replication node3

```

```

postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sync_
nc_priority
FROM pg_stat_replication ;
      pid | username | application_name | client_addr | state |
sync_state | sync_priority
-----+-----+-----+-----+-----+-----+
-----+-----
      26906 | repuser | node3          | 192.168.28.76 |
streaming | quorum   |                | 1
(1 row)

```

□□□□□□□□□□□□□□□□□□□□

```
postgres=# INSERT INTO test_delay(id,create_time)
VALUES(5,now());
--□□□□□□□□
```

□□□□INSERT□□□□□□□□□□□□□□□□□□□□□□□□
□□□'ANY 2□node2□node3□'□□□□□□□□□□□□□□□□
□□□□□□□□□WAL□□□□□□□WAL□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

12.8 高可用

高可用是指系统能够在部分组件发生故障的情况下，仍然能够继续提供服务。PostgreSQL 提供了多种高可用方案，包括主备复制、流复制、WAL 复制、Cascading Replication 等。

12.8.1 主备复制

主备复制是指一个主数据库（Master）与一个或多个备数据库（Slave）保持数据同步。图 12-3 展示了主备复制的架构。

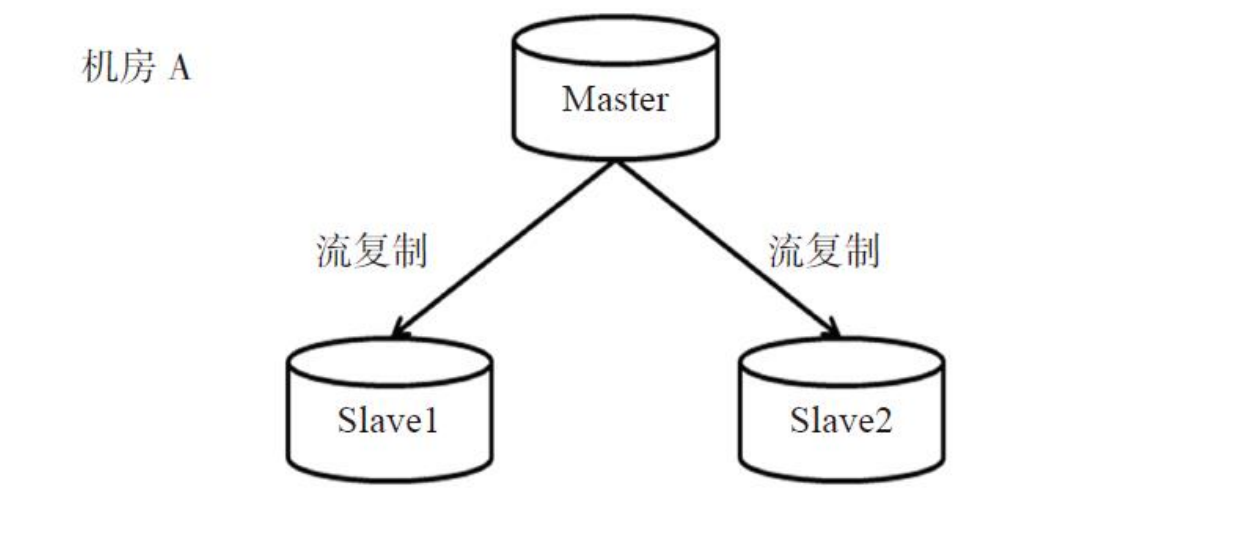


图 12-3 主备复制架构

在主备复制中，Master 数据库负责处理所有的读写操作。Slave1 和 Slave2 数据库只负责接收来自 Master 的数据同步。Slave1 和 Slave2 数据库不能直接接收来自 Master 的写入操作。

机房 A

图 12-4 级联备库

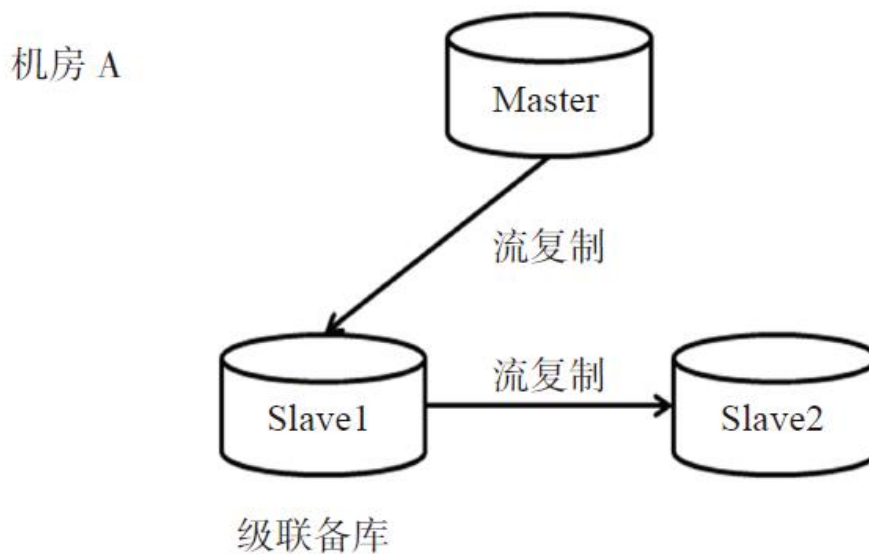


图 12-4 级联备库

图 12-3 展示了 Master 和 Slave2 之间的流复制。Master 和 Slave1 之间的流复制，以及 Slave1 和 Slave2 之间的流复制，都是基于 Master 的 WAL 日志。Master 的 WAL 日志会被复制到 Slave1 的 WAL 日志，而 Slave1 的 WAL 日志会被复制到 Slave2 的 WAL 日志。这种配置称为 cascading standby。

级联备库

· 级联备库 CPU 占用

- [illegible]

机房 A



机房 B



```

Slave1 Master
14 PostgreSQL Slave2

```

12.8.2 五五五五五

12-5

图12-5 物理备份

主 机	主 机 名	IP 地址	操 作 系 统	PostgreSQL 版本
Master	pghost1	192.168.28.74	CentOS6.9	PostgreSQL10
Slave1	pghost2	192.168.28.75	CentOS6.9	PostgreSQL10
Slave2	pghost3	192.168.28.76	CentOS6.9	PostgreSQL10

图12-6 级联备库

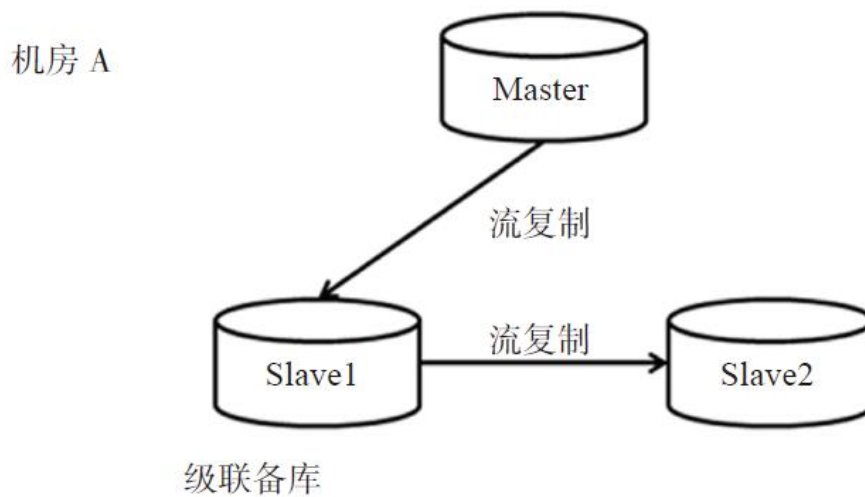


图12-6 级联备库

在Slave1上配置Slave2的备份信息

在Slave1的recovery.conf中配置

```
recovery_target_timeline = 'latest'
standby_mode = on
primary_conninfo = 'host=192.168.28.74 port=1921
user=repuser application_name=slave1'
```

Slave1 Slave2Slave2

```
[postgres@pghost3 pg10]$ pg_basebackup -D
/database/pg10/pg_root -Fp -Xs -v -P -h 192.168.28.74 -p
1921 -U repuser
pg_basebackup: initiating base backup, waiting for
checkpoint to complete
pg_basebackup: checkpoint completed
pg_basebackup: write-ahead log start point: 4/4E000028 on
timeline 7
pg_basebackup: starting background WAL receiver
9424317/9424317 kB (100%), 3/3 tablespaces
pg_basebackup: write-ahead log end point: 4/4E004AA8
pg_basebackup: waiting for background process to finish
streaming ...
pg_basebackup: base backup completed
```

Slave2recovery.conf

```
recovery_target_timeline = 'latest'
standby_mode = on
primary_conninfo = 'host=192.168.28.75 port=1921
user=repuser application_name=slave2'
```

Slave2

```
[postgres@pghost3 pg_root]$ pg_ctl start
```

Slave2

Master pg_stat_replication

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sync_priority
FROM pg_stat_replication ;
 pid | username | application_name | client_addr | state | sync_state | sync_priority
-----+-----+-----+-----+-----+-----+-----
 25041 | repuser | slave1          | 192.168.28.75 | streaming | async | 0
(1 row)
```

Master Slave1 WAL

Slave1 pg_stat_replication

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sync_priority
FROM pg_stat_replication ;
 pid | username | application_name | client_addr | state | sync_state | sync_priority
-----+-----+-----+-----+-----+-----+-----
 5002 | repuser | slave2          | 192.168.28.76 | streaming | async | 0
(1 row)
```

Slave1 Slave2 WAL
Slave1 WAL

Master

```
[postgres@pghost1 ~]$ psql postgres postgres
postgres=# CREATE table t_sr6(id int4);
CREATE TABLE
postgres=# INSERT INTO t_sr6 VALUES (1);
INSERT 0 1
```

Slave1

```
[postgres@pghost2 pg_root]$ psql postgres postgres
postgres=# SELECT * FROM t_sr6;
 id
----
  1
(1 row)
```

Slave1 t_sr6 Slave2

```
[postgres@pghost3 pg_root]$ psql postgres postgres
postgres=# SELECT * FROM t_sr6;
 id
----
  1
(1 row)
```

Slave2

12.9 高可用数据库

PostgreSQL 9.0 高可用数据库解决方案
高可用数据库解决方案“高”是指数据库的高可用性，即数据库在发生故障时能够快速恢复，保证数据的完整性和一致性。

12.9.1 高可用数据库解决方案

高可用数据库解决方案是指通过多种技术手段，保证数据库在发生故障时能够快速恢复，保证数据的完整性和一致性。高可用数据库解决方案通常包括以下几种技术：
1. 主备复制：通过在主数据库和备数据库之间建立复制关系，当主数据库发生故障时，备数据库可以快速接管，保证数据的完整性和一致性。
2. 日志流复制：通过在主数据库和备数据库之间建立日志流复制关系，当主数据库发生故障时，备数据库可以快速接管，保证数据的完整性和一致性。
3. 分布式数据库：通过分布式数据库技术，将数据分布在多个节点上，当某个节点发生故障时，其他节点可以快速接管，保证数据的完整性和一致性。
4. 数据库集群：通过数据库集群技术，将数据库分布在多个节点上，当某个节点发生故障时，其他节点可以快速接管，保证数据的完整性和一致性。
5. 数据库虚拟化：通过数据库虚拟化技术，将数据库分布在多个节点上，当某个节点发生故障时，其他节点可以快速接管，保证数据的完整性和一致性。
6. 数据库备份：通过数据库备份技术，定期备份数据库数据，当数据库发生故障时，可以快速恢复数据，保证数据的完整性和一致性。

高可用数据库解决方案通常包括以下几种技术：
pghost1 和 pghost2 是高可用数据库解决方案中的两个重要组件，它们分别负责主数据库和备数据库的管理。

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sy
nc_priority
FROM pg_stat_replication ;
 pid | username | application_name | client_addr |
state  | sync_state | sync_priority
-----+-----+-----+-----+-----+-----+-----
```

```
--+-----+-----
      24924| repuser | node2          |192.168.28.75|
streaming | async      |          0
(1 row)
```

```

      tbs_his
pgghost1

```

```
[postgres@pgghost1 ~]$ mkdir -p /database/pg10/pg_tbs/tbs_his
```

```

      tbs_his

```

```
postgres=# CREATE TABLESPACE tbs_his OWNER pguser LOCATION
'/database/pg10/pg_tbs/tbs_his';
CREATE TABLESPACE
```

```


```

```
2017-09-12 16:35:11.962 CST,,,16742,,59b79b9f.4166,6,,2017-
09-12 16:32:31 CST,
1/0,0,FATAL,58P01,"directory
"/database/pg10/pg_tbs/tbs_his" does not exist",,"Create
this directory for the tablespace before restarting the
server.",,,"WAL redo at 3/AA17E8D0 for Tablespace/CREATE:
25227 "/database/pg10/pg_tbs/tbs_his""",,,"
2017-09-12 16:35:11.962 CST,,,16740,,59b79b9f.4164,3,,2017-
09-12 16:32:31 CST,,0,LOG,00000,"startup process (PID 16742)
exited with exit code 1",,,,,,,,,,"
2017-09-12 16:35:11.963 CST,,,16740,,59b79b9f.4164,4,,2017-
09-12 16:32:31 CST,,0,LOG,00000,"terminating any other
active server processes",,,,,,,,,,"
2017-09-12 16:35:11.971 CST,,,16740,,59b79b9f.4164,5,,2017-
```

```
09-12 16:32:31 CST,,0,LOG,00000,"database system is shut
down",,,,,,,,,,""
```

```

database/pg10/pg_tbs/tbs_his

```

```
[postgres@pghost2 ~]$ mkdir -p
/database/pg10/pg_tbs/tbs_his
```

```
pghost2
```

```
[postgres@pghost2 ~]$ pg_ctl start
```

```

WAL

```

```
postgres=# SELECT
pid,username,application_name,client_addr,state,sync_state,sy
nc_priority
FROM pg_stat_replication ;
 pid | username | application_name | client_addr |
state | sync_state | sync_priority
-----+-----+-----+-----+-----+-----+-----
 26841| repuser  | node2            | 192.168.28.75 |
streaming | async      |                  | 0
(1 row)
```

数据库的并发控制

数据库的并发控制是指数据库管理系统在并发操作环境下，保证数据库的一致性和完整性。数据库的并发控制是数据库系统的重要组成部分，它涉及到数据库的锁管理、事务管理、死锁检测等方面。

12.9.2 数据库的并发控制

数据库的并发控制是指数据库管理系统在并发操作环境下，保证数据库的一致性和完整性。数据库的并发控制是数据库系统的重要组成部分，它涉及到数据库的锁管理、事务管理、死锁检测等方面。

```
ERROR: canceling statement due to conflict with recovery
DETAIL: User query might have needed to see row versions
that must be removed.
```

数据库的并发控制是指数据库管理系统在并发操作环境下，保证数据库的一致性和完整性。数据库的并发控制是数据库系统的重要组成部分，它涉及到数据库的锁管理、事务管理、死锁检测等方面。

PostgreSQL MVCC 数据库的并发控制

PostgreSQL 数据库的并发控制是通过 MVCC (Multi-Version Concurrency Control) 实现的。MVCC 是一种数据库并发控制技术，它允许数据库在并发操作环境下，保证数据库的一致性和完整性。

PostgreSQL 数据库的并发控制是通过 MVCC (Multi-Version Concurrency Control) 实现的。MVCC 是一种数据库并发控制技术，它允许数据库在并发操作环境下，保证数据库的一致性和完整性。

PostgreSQL 数据库的并发控制是通过 MVCC (Multi-Version Concurrency Control) 实现的。MVCC 是一种数据库并发控制技术，它允许数据库在并发操作环境下，保证数据库的一致性和完整性。

·max_standby_streaming_delay
30SQLWAL
3030WALWAL
-1WALWAL

·hot_standby_feedback
on
off

pghost1pghost2
postgresql.conf

```
max_standby_streaming_delay = 10s
```

max_standby_streaming_delay10
reload

```
[postgres@pghost2 pg_root]$ pg_ctl reload
server signaled
```

update_per2.sql

```
\set v_id random(1,1000000)
update test_per2 set flag='1' where id=:v_id;
```

pghost1 pgbench 120

```
pgbench -c 8 -T 120 -d postgres -U postgres -n N -M
prepared -f update_per2.sql> update_8.out 2>&1 &
```

```
postgres=# \timing
Timing is on.
postgres=# SELECT pg_sleep(15),count(*) FROM test_per2;
ERROR: canceling statement due to conflict with recovery
DETAIL: User query might have needed to see row versions
that must be removed.
Time: 10433.102 ms (00:10.433)
```

test_per2 pg_sleep

max_standby_streaming_delay
10

max_standby_streaming_delay
10 WAL SQL
10 -1
60

```
max_standby_streaming_delay = 60s  
hot_standby_feedback = off
```

hot_standby_feedback off
reload

```
[postgres@pghost2 pg_root]$ pg_ctl reload  
server signaled
```

pgbench

```
postgres=# SELECT pg_sleep(15),count(*) FROM test_per2;  
pg_sleep | count  
-----+-----  
          | 10000000  
(1 row)
```

```
Time: 15327.394 ms (00:15.327)
```

15

配置hot_standby_feedback

hot_standby_feedback 配置 on
配置 hot_standby_feedback 为 on
配置 max_standby_streaming_delay 为 10s

```
hot_standby_feedback = on
max_standby_streaming_delay = 10s
```

配置 hot_standby_feedback 为 on
配置 max_standby_streaming_delay 为 10s
配置 reload

```
[postgres@pghost2 pg_root]$ pg_ctl reload
server signaled
```

配置 pgbench 配置

```
postgres=# SELECT pg_sleep(15),count(*) FROM test_per2;
 pg_sleep | count
-----+-----
          | 10000000
(1 row)
```

Time: 15349.958 ms (00:15.350)

配置 15s

max_standby_streaming_delay=1
hot_standby_feedback

12.9.3 WAL

FATAL,XX000,"could not receive data from WAL stream: ERROR: requested WAL segment 000000010000000100000022 has already been removed"

WAL
000000010000000100000022
WAL
WAL
wal_keep_segment
WAL WAL

データベースのサイズは1TB、バックアップは600GB、
バックアップは毎日1回、バックアップは毎日1回、
バックアップは毎日1回、バックアップは毎日1回、

バックアップは毎日1回、バックアップは毎日1回、
バックアップは毎日1回、バックアップは毎日1回、
postgresql.confバックアップは毎日1回、

```
wal_keep_segments = 1  
archive_mode = on  
archive_command = 'cp %p /archive_dir/pg10/%f'
```

wal_keep_segmentsは1、pg_wal
WALバックアップは毎日1回、
WALバックアップは毎日1回、
reloadバックアップは毎日1回、

```
[postgres@pghost1 pg_root]$ pg_ctl reload  
server signaled
```

wal_keep_segmentsは1、
pg_walは100、WALバックアップは毎日1回、
checkpointはpg_walバックアップは毎日1回、
バックアップは毎日1回、

```
postgres=# CHECKPOINT;
CHECKPOINT
```

```
pg_walWALWAL
WAL
```

```
[postgres@pghost2 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

```
update_per2.sql
```

```
\set v_id random(1,1000000)
```

```
update test_per2 set flag='1' where id=:v_id;
```

```
pgbench
```

```
pgbench -c 8 -T 120 -d postgres -U postgres -n N -M prepared
-f update_per2.sql > update_8.out 2>&1 &
[1] 23803
```

```
pgbench
pg_switch_walcheckpoint
```

□□□□□

```
postgres=# SELECT pg_switch_wal();
           pg_switch_wal
-----
          4/32983D58
(1 row)
...□□□□
```

```
postgres=# checkpoint;
CHECKPOINT
...□□□□
```

□□□□□□□□□□pg_switch_wal□□□□□□□□□□
WAL□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
[postgres@pghost1 pg_root]$ ll /archive_dir/pg10/
total 352M
-rw----- 1 postgres postgres 16M Sep 14 21:34
0000000700000000400000023
-rw----- 1 postgres postgres 16M Sep 14 21:34
0000000700000000400000024
□□
```

□□WAL□□□□/archive_dir/pg10/□□□□□□□□
WAL□□□□□□

□□□□□□□□□□□□□□

```
[postgres@pghost2 ~]$ pg_ctl start
server started
```

WAL
[redacted]

```
[postgres@pghost1 pg_root]$ scp /archive_dir/pg10/*  
postgres@pghost2:/archive_dir/pg10  
postgres@pghost2's password:
```

recovery.conf

```
restore_command = 'cp /archive_dir/pg10/%f %p'
```

restore_command shell
WAL WAL

```
[postgres@pghost2 pg_root]$ pg_ctl restart
```

[redacted]

```
2017-09-14 21:55:16.904 CST,,,25224,,59ba8a44.6288,1,,2017-  
09-14 21:55:16 CST,,0,LOG,00000,"ending log output to  
stderr",,"Future log output will go to log destination  
""csvlog"".",,,,,,,,,""  
2017-09-14 21:55:16.911 CST,,,25226,,59ba8a44.628a,1,,2017-  
09-14 21:55:16 CST,,0,LOG,00000,"database system was shut  
down in recovery at 2017-09-14 21:55:11 CST",,,,,,,,,""  
2017-09-14 21:55:16.913 CST,,,25226,,59ba8a44.628a,2,,2017-  
09-14 21:55:16 CST,,0,LOG,00000,"entering standby  
mode",,,,,,,,,""  
2017-09-14 21:55:16.933 CST,,,25226,,59ba8a44.628a,3,,2017-  
09-14 21:55:16 CST,1/0,0,LOG,00000,"redo starts at
```

```
4/1A00D8A8",,,,,,,,,,""
2017-09-14 21:55:18.052 CST,,,25226,,59ba8a44.628a,4,,2017-
09-14 21:55:16 CST,1/0,0,LOG,00000,"restored log file
""0000000700000000400000023"" from archive",,,,,,,,,,""
2017-09-14 21:55:18.526 CST,,,25226,,59ba8a44.628a,5,,2017-
09-14 21:55:16 CST,1/0,0,LOG,00000,"consistent recovery
state reached at 4/23D6EA98",,,,,,,,,,""
2017-09-14 21:55:18.527 CST,,,25224,,59ba8a44.6288,2,,2017-
09-14 21:55:16 CST,,0,LOG,00000,"database system is ready to
accept read only connections",,,,,,,,,,""
2017-09-14 21:55:18.552 CST,,,25226,,59ba8a44.628a,6,,2017-
09-14 21:55:16 CST,1/0,0,LOG,00000,"restored log file
""0000000700000000400000024"" from archive",,,,,,,,,,""
2017-09-14 21:55:18.685 CST,,,25226,,59ba8a44.628a,7,,2017-
09-14 21:55:16 CST,1/0,0,LOG,00000,"restored log file
""0000000700000000400000025"" from archive",,,,,,,,,,""
...
2017-09-14 21:55:22.192 CST,,,25264,,59ba8a4a.62b0,1,,2017-
09-14 21:55:22 CST,,0,LOG,00000,"started streaming WAL from
primary at 4/38000000 on timeline 7",,,,,,,,,,""
```

```

WAL
0000000700000000400000023
WAL"started
streaming WAL"
```

```

WAL
WAL
```

```


```

```

wal_keep_segments
$PGDATA/pg_walWAL
```


pg_wal WAL 目录
pg_wal 目录

pg_wal 目录
pg_wal 目录
pg_wal 目录
pg_wal 目录

pg_wal Replication slots
pg_wal
pg_wal
pg_wal
pg_wal
pg_wal
pg_wal
pg_wal

postgresql.conf
pg_wal

```
max_replication_slots = 1  
wal_keep_segments = 1
```

pg_wal

pg_wal

```

postgres=# SELECT * FROM
pg_create_physical_replication_slot('phy_slot1');
   slot_name | lsn
-----+-----
   phy_slot1 |
(1 row)

```

pg_replication_slots
 pg_replication_slots 테이블은 PostgreSQL 10 버전에서 도입된 테이블로, 물리적 복제 슬롯의 정보를 저장한다.

```

postgres=# SELECT
slot_name,plugin,slot_type,database,active,active_pid,xmin
FROM pg_replication_slots ;
   slot_name | plugin | slot_type | database | active |
active_pid | xmin
-----+-----+-----+-----+-----+-----
   phy_slot1 |        | physical  |          | f      |
          |
(1 row)

```

- slot_name: 슬롯의 이름
- plugin: 슬롯의 플러그인 이름
- slot_type: 슬롯의 유형 (physical, logical)
- database: 슬롯이 복제하는 데이터베이스 이름
- active: 슬롯이 활성 상태인지를 나타내는 비트

·active_pid

·xmin

recovery.conf
primary_slot_name

```
recovery_target_timeline = 'latest'  
standby_mode = on  
primary_conninfo = 'host=192.168.28.74 port=1921  
user=repuser application_name=slavel'  
primary_slot_name = 'phy_slot1'
```

```
[postgres@pghost2 pg_root]$ pg_ctl restart -m fast  
waiting for server to shut down.... done  
server stopped
```

pg_replication_slots

```
slot_name,plugin,slot_type,database,active,active_pid,xmin  
FROM pg_replication_slots ;  
slot_name | plugin | slot_type | database | active |  
active_pid | xmin  
-----+-----+-----+-----+-----+-----  
phy_slot1 |      | physical |      | t      |  
23833 |  
(1 row)
```

active transactions 23833

```
[postgres@pghost1 pg_root]$ ps -ef | grep 23833
postgres 23833 19503  0 14:40 ?        00:00:00 postgres:
wal sender process repuser 192.168.28.75(33141) streaming
4/52044950
```

23833 WAL

WAL

```
[postgres@pghost2 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

pgbench

```
pgbench -c 8 -T 120 -d postgres -U postgres -n N -M prepared
-f update_per2.sql > update_8.out 2>&1 &
```

pg_switch_wal
checkpoint

pg_ctl start

```
[postgres@pgghost2 ~]$ pg_ctl start
```

```

    WAL
$PGDATA/pg_wal
WAL

```

12.10 概要

PostgreSQL10は、8つの
PostgreSQL10のLogical
Replicationの2ndquadrant
を実装する。

この機能は、
PostgreSQL10のLogical
Replicationの2ndquadrant
を実装する。

- ・PostgreSQL10のLogical Replicationの2ndquadrantを実装する。
- ・PostgreSQL10のLogical Replicationの2ndquadrantを実装する。
- ・PostgreSQL10のLogical Replicationの2ndquadrantを実装する。
- ・PostgreSQL10のLogical Replicationの2ndquadrantを実装する。

PostgreSQL10のWALは、
WALのLogical decoding
WALのLogical decoding
WALのLogical decoding
WALのLogical decoding
WALのLogical decoding

12.10.1 配置

配置logical decoding需要配置以下参数：
wal_level=logical
max_replication_slots=1

配置wal_level=logical需要配置以下参数：
max_replication_slots=1

```
wal_level = logical  
max_replication_slots = 8
```

wal_level=WAL
replica_logical=12.1.1
配置logical decoding

max_replication_slots=1
配置logical decoding

pghost1

```
postgres=# SELECT  
pg_create_logical_replication_slot('logical_slot1','test_dec  
oding');  
 slot_name | lsn  
-----+-----  
 logical_slot1 | 4/85004210  
(1 row)  
pg_replication_slots  
postgres=# SELECT
```

```

slot_name,plugin,slot_type,database,active,restart_lsn
      FROM pg_replication_slots;
 slot_name | plugin | slot_type | database |
active | restart_lsn
-----+-----+-----+-----+---
logical_slot1 | test_decoding | logical | postgres | f
| 4/850041D8
phy_slot1 | | physical | | t
| 4/85016670
(2 rows)

```

logical_slot1
 pg_h_slot1 12.9.3
 WAL

logical_slot1

```

postgres=# SELECT * FROM
pg_logical_slot_get_changes('logical_slot1',null,null);
 lsn | xid | data
-----+-----+-----
(0 rows)

```

pg_logical_slot_get_changes

DDL

DDL DDL

[illegible]

```
postgres=# INSERT INTO t_logical VALUES (1);
INSERT 0 1
postgres=# SELECT * FROM
pg_logical_slot_get_changes('logical_slot1',null,null);
      lsn      |      xid      |      data
-----+-----+-----
-----+-----+-----
```

```

4/850AB898 | 42976848 | BEGIN 42976848
4/850AB898 | 42976848 | table public.t_logical: INSERT:
id[integer]:1
4/850AB908 | 42976848 | COMMIT 42976848
(3 rows)

```

pg_logical_slot_get_changes



pg_logical_slot_get_changes

pg_logical_slot_peek_changes

pg_logical_slot_get_changes

pg_logical_slot_peek_changes

pg_logical_slot_get_changes

pg_logical_slot_peek_changes

pg_logical_slot_peek_changes

pg_recvlogical

pg_recvlogical

pg_recvlogical

pg_recvlogical

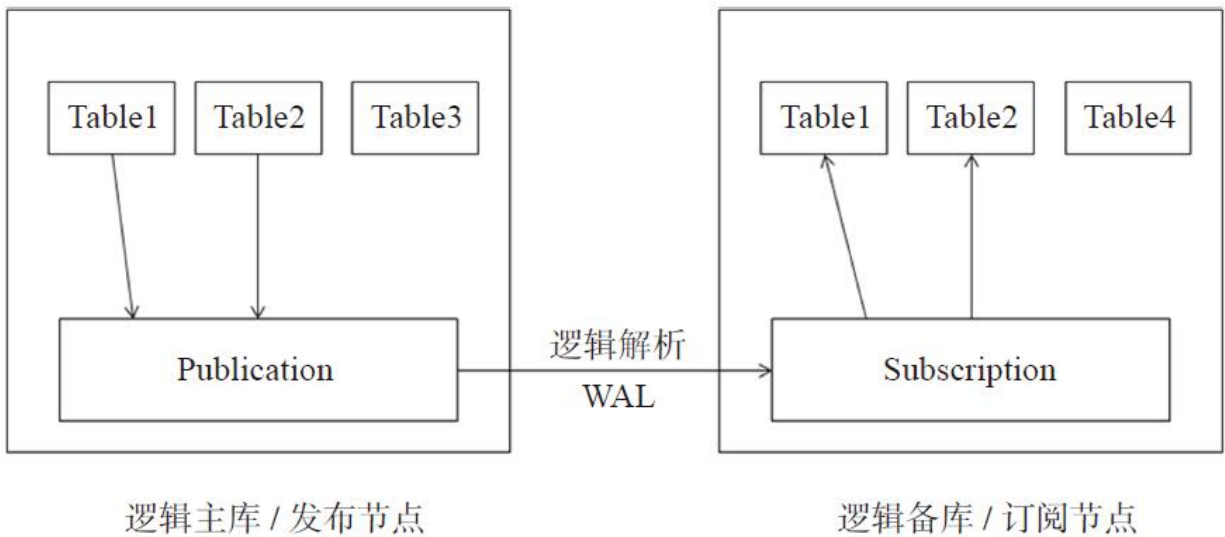
```

postgres=# INSERT INTO t_logical VALUES (3);
INSERT 0 1

```

pg_recvlogical

logical_slot1



12-7 逻辑复制

逻辑复制是 PostgreSQL 9.3 版本引入的新特性，它允许将一个数据库中的表复制到另一个数据库中。在逻辑复制中，主库通过 **Publication** 组件将表的数据复制到备库的 **Subscription** 组件。主库的 **Publication** 组件会记录表 **table1**、**table2** 等的变更，并将这些变更通过 **WAL** 发送给备库的 **Subscription** 组件。备库的 **Subscription** 组件会根据接收到的 **WAL** 记录，在备库中重新创建表 **table1**、**table2** 等。

逻辑复制的架构如下图所示。主库的 **Publication** 组件会记录表 **table1**、**table2** 等的变更，并将这些变更通过 **WAL** 发送给备库的 **Subscription** 组件。备库的 **Subscription** 组件会根据接收到的 **WAL** 记录，在备库中重新创建表 **table1**、**table2** 等。

逻辑复制的架构如下图所示。主库的 **Publication** 组件会记录表 **table1**、**table2** 等的变更，并将这些变更通过 **WAL** 发送给备库的 **Subscription** 组件。备库的 **Subscription** 组件会根据接收到的 **WAL** 记录，在备库中重新创建表 **table1**、**table2** 等。

逻辑复制的架构如下图所示。主库的 **Publication** 组件会记录表 **table1**、**table2** 等的变更，并将这些变更通过 **WAL** 发送给备库的 **Subscription** 组件。备库的 **Subscription** 组件会根据接收到的 **WAL** 记录，在备库中重新创建表 **table1**、**table2** 等。

replica
identity 删除/UPDATE
full 更新/DELETE

Subscription
Subscription
DDL
DDL WAL

12.10.3

12-6

12-6

角 色	主 机 名	IP	端 口	库 名	用 户 名	版 本
发布节点	pghost1	192.168.28.74	1921	mydb	pguser	PostgreSQL10
订阅节点	pghost3	192.168.28.76	1923	des	pguser	PostgreSQL10

postgresql.conf

```
wal_level = logical  
max_replication_slots = 8
```

```
max_wal_senders = 10
```

```
·wal_levellogical
·max_replication_slots
·max_wal_senders
WAL
max_replication_slots
postgresql.conf
```

```
max_replication_slots = 8
max_logical_replication_workers = 8      # taken from
max_worker_processes
```

```
·max_replication_slots
```

```
·max_logical_replication_workers
4
```

```
max_logical_replication_workers
max_worker_processes
max_worker_processes
```

REPLICATION

```
CREATE USER logical_user
  REPLICATION
  LOGIN
  CONNECTION LIMIT 8
  ENCRYPTED PASSWORD 'logical_user';
```

```
logical_user REPLICATION
logical_user SUPERUSER
logical_user logical_user
```

```
[postgres@pghost1 ~]$ psql mydb pguser
psql (10.0)
Type "help" for help.
```

```
mydb=> CREATE TABLE t_lr1(id int4,name text);
CREATE TABLE
mydb=> INSERT INTO t_lr1 VALUES (1,'a');
INSERT 0 1
```

```
t_lr1
```

```
mydb=> CREATE PUBLICATION pub1 FOR TABLE t_lr1;
CREATE PUBLICATION
```

CREATE PUBLICATION

```
CREATE PUBLICATION name
  [ FOR TABLE [ ONLY ] table_name [ * ] [, ...]
    | FOR ALL TABLES ]
  [ WITH ( publication_parameter [= value] [, ... ] ) ]
```

·name

·FOR TABLE
指定するテーブルをこのパブリケーションに含める。この場合、パブリケーションは、指定されたテーブルのみに適用される。パブリケーションは、指定されたテーブルのみに適用される。

·FOR ALL TABLES
このパブリケーションは、データベース内のすべてのテーブルに適用される。PostgreSQL 10.0以降では、このオプションは省略可能である。PostgreSQL 10.0以前では、このオプションは必須である。

pg_publication

```
mydb=> SELECT * FROM pg_publication;
      pubname | pubowner | puballtables | pubinsert |
pubupdate | pubdelete
-----+-----+-----+-----+-----
--+-+-----+
      pub1   |    16391 | f           | t         |
| t
(1 row)
```

- pubname 发布者名称
- pubowner 发布者pg_user 用户名
usesysid 用户ID
- puballtables 发布的所有表名 t 表名
发布的所有表名
- pubinsert t 表名 INSERT 语句
- pubupdate t 表名 UPDATE 语句
- pubdelete t 表名 DELETE 语句

发布者 t_lr1 表名

```
[des@pghost3 ~]$ psql des pguser
psql (10.0)
Type "help" for help.
```

```
des=> CREATE TABLE t_lr1(id int4,name text);
CREATE TABLE
```

发布者名称

```
CREATE SUBSCRIPTION subscription_name
CONNECTION 'conninfo'
PUBLICATION publication_name [, ...]
[ WITH ( subscription_parameter [= value] [, ... ] ) ]
```

·subscription_name

·CONNECTION host
port dbname user password
 ~/.pgpass

·PUBLICATION

·WITH subscription_parameter[=value]
[...]
copy_data boolean
create_slot boolean enabled boolean
slot_name string

~/.pgpass

192.168.28.74:1921:mydb:logical_user:logical_user

~/.pgpass

[des@pgghost3 ~]\$ chmod 0600 .pgpass

pg_hba.conf

```
[des@pgghost3 ~]$ psql des postgres
psql (10.0)
Type "help" for help.
```

```
des=# CREATE SUBSCRIPTION sub1 CONNECTION
'host=192.168.28.74 port=1921
      dbname=mydb user=logical_user' PUBLICATION pub1;
NOTICE:  created replication slot "sub1" on publisher
CREATE SUBSCRIPTION
```

```

sub1
```

```
mydb=> SELECT
slot_name,plugin,slot_type,database,active,restart_lsn
      FROM pg_replication_slots where slot_name='sub1';
      slot_name | plugin | slot_type | database | active |
restart_lsn
-----+-----+-----+-----+-----+--
      sub1      | pgoutput | logical   | mydb     | t      |
4/9793A6F0
(1 row)
```

```

plugin
pgoutput
```

```

pg_subscription

```

```
des=# SELECT * FROM pg_subscription;
-[ RECORD 1 ]-----

```

subdbid		16387
subname		sub1
subowner		10
subenabled		t
subconninfo		host=192.168.28.74 port=1921 dbname=mydb
user=logical_user		
subslotname		sub1
subsynccommit		off
subpublications		{pub1}

·subdbid OID pg_database.oid

·subname

·subowner

·subenabled

·subconninfo

·subslotname

·subpublications

t_lr1

```
des=> SELECT * FROM t_lr1;
id | name
----+-----
(0 rows)
```

pguser t_lr1
pg

```
2017-10-01 14:39:45.795 CST,,,16650,,59d08db1.410a,1,,2017-
10-01 14:39:45 CST,4/47,0,LOG,00000,"logical replication
table synchronization worker for subscription "sub1",
table "t_lr1" has started",,,,,,""
2017-10-01 14:39:45.875 CST,,,16650,,59d08db1.410a,2,,2017-
10-01 14:39:45 CST,4/50,0,ERROR,XX000,"could not start
initial contents copy for table "pguser.t_lr1": ERROR:
permission denied for schema pguser",,,,,,""
2017-10-01 14:39:45.875 CST,,,6054,,59cf4455.17a6,286,,2017-
09-30 15:14:29 CST,,0,LOG,00000,"worker process: logical
replication worker for subscription 16396 sync 16389 (PID
16650) exited with exit code 1",,,,,,""
```

pguser sub1 t_lr1
pguser

logical_user
logical_user

```
mydb=> GRANT USAGE ON SCHEMA pguser TO logical_user;
GRANT
mydb=> GRANT SELECT ON t_lr1 TO logical_user;
GRANT
```

pguser logical_user
t_lr1 SELECT logical_user

```
2017-10-01 14:45:18.893 CST,,,16755,,59d08efe.4173,1,,2017-
10-01 14:45:18 CST,4/400,0,LOG,00000,"logical replication
table synchronization worker for subscription ""sub1"",
table ""t_lr1"" has started",,,,,,,,,,""
2017-10-01 14:45:18.951 CST,,,16755,,59d08efe.4173,2,,2017-
10-01 14:45:18 CST,4/404,0,LOG,00000,"logical replication
table synchronization worker for subscription ""sub1"",
table ""t_lr1"" has finished",,,,,,,,,,""
```

```
sub1t_lr1

```

```
des=> SELECT * from t_lr1 ;
      id | name
-----+-----
       1 | a
(1 row)
```

```
t_lr1
WAL

```

```
postgres: wal sender process logical_user
192.168.28.76(47464) idle
```

```
WAL
```

```
postgres: bgworker: logical replication worker for
subscription 16396
```

데이터베이스에서 DML을 실행할 때

12.10.4 데이터 DML

데이터베이스에서 INSERT/UPDATE/DELETE를 실행할 때

데이터베이스에서

```
mydb=> INSERT INTO t_lr1 VALUES (2,'b');
INSERT 0 1
```

데이터베이스에서

```
des=> SELECT * FROM t_lr1 WHERE id=2;
      id | name
-----+-----
      2  | b
(1 row)
```

데이터베이스에서

```
mydb=> UPDATE t_lr1 SET name='bb' WHERE id=2;
ERROR:  cannot update table "t_lr1" because it does not have
        replica identity and publishes updates
```

HINT: To enable updating the table, set REPLICA IDENTITY using ALTER TABLE.

```

t_lr1 replica
identity 12.10.2
UPDATE/DELETE
replica identity
t_lr1

```

t_lr1

```
mydb=> ALTER TABLE t_lr1 ADD PRIMARY KEY(id);
ALTER TABLE
```

t_lr1

```
des=> ALTER TABLE t_lr1 ADD PRIMARY KEY(id);
ALTER TABLE
```

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □

```
mydb=> UPDATE t_lr1 SET name='bb' WHERE id=2;
UPDATE 1
```

更新表 t_lr1 的 id=2 的记录 UPDATE

UPDATE t_lr1 SET name='bb' WHERE id=2;

```
des=> SELECT * FROM t_lr1 WHERE id=2;
      id | name
-----+-----
      2 | bb
(1 row)
```

更新表 t_lr1 的 id=2 的记录 UPDATE

DELETE FROM t_lr1 WHERE id=2;

删除 ID 为 2 的记录 DELETE

```
mydb=> DELETE FROM t_lr1 WHERE id=2;
DELETE 1
```

删除表 t_lr1 的 id=2 的记录 DELETE

```
des=> SELECT * FROM t_lr1 WHERE id=2;
      id | name
-----+-----
(0 rows)
```

删除表 t_lr1 的 id=2 的记录 DELETE

删除

```

t_lr1
t_lr1

```

12.10.5

```


```

```

t_big1000

```

```

mydb=> CREATE TABLE t_big(id int4 primary key,
create_time timestamp(0) without time zone default
clock_timestamp(),
name character varying(32));)
CREATE TABLE

mydb=> INSERT INTO t_big(id,name)
SELECT n,n*random()*10000 FROM generate_series(1,10000000)
n;
INSERT 0 10000000

```

```

t_bigSELECT
logical_user

```

```

mydb=> GRANT SELECT ON t_big TO logical_user;
GRANT

```

CREATE PUBLICATION pub1 FOR TABLE t_big;

```
mydb=> ALTER PUBLICATION pub1 ADD TABLE t_big;
ALTER PUBLICATION
```

ALTER PUBLICATION pub1 ADD TABLE t_big;

```
mydb=> \dRp+ pub1
          Publication pub1
   Owner  | All tables | Inserts | Updates | Deletes
-----+-----+-----+-----+-----
   pguser | f          | t       | t       | t
Tables:
    "pguser.t_big"
    "pguser.t_lr1"
```

CREATE PUBLICATION pub1 FOR TABLE t_big;

CREATE PUBLICATION pub1 FOR TABLE t_big;

```
mydb=> SELECT * FROM pg_publication_tables ;
 pubname | schemaname | tablename
-----+-----+-----
   pub1  | pguser    | t_lr1
   pub1  | pguser    | t_big
(2 rows)
```

CREATE PUBLICATION pub1 FOR TABLE t_big;

```
des=> CREATE TABLE t_big(id int4 primary key,  
create_time timestamp(0) without time zone default  
clock_timestamp(),  
name character varying(32));  
CREATE TABLE
```

```
    t_big  
    t_big  
    t_big
```

```
des=# ALTER SUBSCRIPTION sub1 REFRESH PUBLICATION ;  
ALTER SUBSCRIPTION
```

```
    t_big  
COPY 33  
t_big
```

```
des=# SELECT COUNT(*) FROM pguser.t_big;  
count  
-----  
10000000  
(1 row)
```

```
    t_big  
  
    t_big  
t_big pub1
```

```
mydb=> ALTER PUBLICATION pub1 DROP TABLE t_big;  
ALTER PUBLICATION
```

```
des=# ALTER SUBSCRIPTION sub1 ENABLE ;
ALTER SUBSCRIPTION
```

```
pg_subscription subenabled
t
```

```
des=# SELECT subname,subenabled,subpublications FROM
pg_subscription;
  subname | subenabled | subpublications
-----+-----+-----
    sub1  | t          | {pub1}
(1 row)
```

12.10.7 复制配置

复制配置包括以下参数：

复制配置参数：

- `wal_level` 参数设置为 `logical`
- `replication` 参数设置为 `yes`
- `UPDATE/DELETE` 操作需要设置 `allow_unaligned_replication` 参数为 `yes`

·INSERT UPDATE DELETE
DML

·

·

·

·pg_hba.conf

·DDL
DDL

·

·

·

·

·

创建表INSERT数据

创建表t_per1并授予logical_user权限

```
mydb=> CREATE TABLE t_per1 (id int4,name text,  
create_time timestamp(0) without time zone DEFAULT '2000-01-  
01 00:00:00');  
CREATE TABLE
```

```
mydb=> GRANT SELECT ON t_per1 TO logical_user;  
GRANT
```

创建发布pub1并添加表t_per1

```
mydb=> ALTER PUBLICATION pub1 ADD TABLE t_per1 ;  
ALTER PUBLICATION
```

创建表t_per1

```
des=> CREATE TABLE t_per1 (id int4,name text,  
create_time timestamp(0) without time zone DEFAULT '2000-01-  
01 00:00:00');  
CREATE TABLE
```

创建订阅sub1

```
des=# ALTER SUBSCRIPTION sub1 REFRESH PUBLICATION ;  
ALTER SUBSCRIPTION
```

pgbench insert_t.sql

```
\set v_id random(1,1000000)

INSERT INTO t_per1(id,name) VALUES (:v_id,:v_id||'a');
```

pgbench INSERT

```
pgbench -c 8 -T 120 -d mydb -U pguser -n N -M prepared -f
insert_t.sql > insert_t.out 2>&1 &
```

INSERT SQL

```
mydb=# SELECT pid,username,state,
           pg_wal_lsn_diff(pg_current_wal_lsn(),replay_lsn)
           replay_dely
           FROM pg_stat_replication WHERE
application_name='sub1';
   pid |   username   | state   | replay_dely
-----+-----+-----+-----
   457 | logical_user | streaming |          11936
(1 row)
```

replay_dely WAL
SQL WAL 11MB

UPDATE

创建表t_per2 1000行

```
mydb=> CREATE TABLE t_per2 (id int4 primary key,  
name text,  
create_time timestamp(0) without time zone DEFAULT '2000-01-  
01 00:00:00');  
CREATE TABLE  
mydb=> INSERT INTO t_per2(id) SELECT  
generate_series(1,1000000);  
INSERT 0 1000000
```

授予权限

```
mydb=> GRANT SELECT ON t_per2 TO logical_user ;  
GRANT  
mydb=> ALTER PUBLICATION pub1 ADD TABLE t_per2;  
ALTER PUBLICATION
```

刷新订阅

```
des=> CREATE TABLE t_per2 (id int4 primary key,  
name text,  
create_time timestamp(0) without time zone DEFAULT '2000-01-  
01 00:00:00');  
CREATE TABLE  
  
des=# ALTER SUBSCRIPTION sub1 REFRESH PUBLICATION ;  
ALTER SUBSCRIPTION
```

40行数据

update_t.sql

```
\set v_id random(1,1000000)
```

```
update t_per2 set create_time=clock_timestamp() where  
id=:v_id;
```

pgbench UPDATE

```
pgbench -c 8 -T 120 -d mydb -U pguser -n N -M prepared -f  
update_t.sql > update_t.out 2>&1 &
```

UPDATE SQL

```
mydb=# SELECT pid,username,state,  
              pg_wal_lsn_diff(pg_current_wal_lsn(),replay_lsn)  
replay_dely  
              FROM pg_stat_replication WHERE  
application_name='sub1';  
   pid |   username   |   state   | replay_dely  
-----+-----+-----+-----  
   457 | logical_user | streaming |    12352
```

SQL replay_dely 12MB

12.11 □□□□

```

graph LR
    DB[(PostgreSQL)] --- App[Java]
  
```

13

13.1 数据库安全

数据库安全是指保护数据库中的数据不被非法访问、篡改、破坏或泄露。数据库安全包括物理安全、逻辑安全、操作安全、应用安全等方面。物理安全是指保护数据库的物理设备（如硬盘、服务器）免受自然灾害、火灾、盗窃等威胁。逻辑安全是指保护数据库的逻辑结构（如表、视图、索引）免受非法访问、篡改、破坏等威胁。操作安全是指保护数据库的操作过程（如备份、恢复、迁移）免受非法访问、篡改、破坏等威胁。应用安全是指保护数据库的应用程序（如Web应用、移动应用）免受非法访问、篡改、破坏等威胁。

数据库安全的主要威胁包括：非法访问、篡改、破坏、泄露、病毒、木马、蠕虫、恶意软件等。数据库安全的主要措施包括：访问控制、身份认证、授权管理、数据加密、备份恢复、灾难恢复、安全审计、漏洞扫描、渗透测试等。

数据库安全是一个持续的过程，需要不断地进行安全评估、安全加固、安全监测、安全响应等工作。

- 数据库安全审计

- 数据库安全加固“三原则”
TRUNCATE
ALTER
TABLE DROP COLUMN

- 数据库安全加固“三原则”

- 数据库安全加固“三原则”
BUG
WHERE
UPDATE
DELETE

- 数据库安全加固“三原则”
SQL注入攻击

数据库安全是一个持续的过程，需要不断地进行安全评估、安全加固、安全监测、安全响应等工作。

1. 在系统启动时，CPU 会执行 BIOS 中的 POST 程序，检查硬件设备是否正常工作。如果检测到硬件故障，CPU 会发出报警声，并显示错误信息。

2. CPU 会根据 BIOS 中的设置，加载操作系统。操作系统启动后，CPU 会执行操作系统中的各种程序，完成系统的正常运行。

3. CPU 会根据操作系统的调度，执行各种任务。CPU 会根据任务的优先级，分配时间片，确保系统的高效运行。

PostgreSQLバックアップと復元の方法について
バックアップと復元の方法についてSQLバックアップと復元
バックアップと復元の方法についてバックアップと復元の方法について
バックアップと復元の方法についてPostgreSQLバックアップ

・pg_dumpとpg_dumpallでバックアップとSQLバックアップ
バックアップ

・バックアップと復元

・バックアップと復元PITR

バックアップと復元バックアップと復元

13.2

PostgreSQL WAL
PostgreSQL WAL
PostgreSQL WAL
WAL
N
 $N = 2 \times \text{checkpoint_segment} + 1$
PostgreSQL WAL
256 WAL 1 2 3
WAL WAL
1 WAL
WAL

The diagram illustrates a sequence of 100 blocks arranged in 10 rows of 10. The blocks are labeled with 'WAL' and 'DBA' in a repeating pattern. The first row starts with 'WAL' at the 8th block. The second row starts with 'WAL' at the 1st block. The third row starts with 'DBA' at the 7th block. The fourth row starts with 'DBA' at the 8th block. The fifth row starts with 'DBA' at the 1st block. The sixth row starts with 'DBA' at the 8th block. The seventh row starts with 'DBA' at the 1st block. The eighth row starts with 'DBA' at the 8th block. The ninth row starts with 'DBA' at the 1st block. The tenth row starts with 'DBA' at the 8th block.

pgdata 目录结构如下

13.2.1 WAL 目录

1. 目录结构

pgdata 目录结构如下

```
data backups
scripts archive_wals
```

```
[root@pghost1 ~]$ mkdir -p
/pgdata/10/{data,backups,scripts,archive_wals}
[root@pghost1 ~]$ chown -R postgres.postgres /pgdata/10
```

data backups
scripts
archive_wals
NFS postgres

2. wal_level

wal_level minimal replica
logical minimal replica logical WAL
WAL minimal
WAL WAL

minimal archive_mode WAL
wal_level replica

```
mydb=# ALTER SYSTEM SET wal_level = 'replica';  
ALTER SYSTEM
```

3. archive_mode

archive_mode on off always
off on

```
mydb=# ALTER SYSTEM SET archive_mode = 'on';  
ALTER SYSTEM
```

4. archive_command

archive_command
shell shell
archive_command shell "%p"
WAL "%f"
WAL

archive_command
archive_command='cp%p/pgdata/10/archi

ve_wals/%f'

```
wal_levelarchive_mode
archive_command
reload
archive_command
WAL
pg_walpg_wal
wal_level
archive_mode
archive_command
/bin/true
archive_commandreload
```

```
WAL
gzipbzip2lz4
archive_commandpg_wallz4
WAL
/pgdata/10/archive_wals/
```

```
mydb=# ALTER SYSTEM SET archive_command = '/usr/bin/lz4 -q -
z %p /pgdata/10/archive_wals/%f.lz4';
ALTER SYSTEM
mydb=# SELECT pg_reload_conf();
pg_reload_conf
-----
```

```

t
(1 row)
mydb=# show archive_command ;
        archive_command
-----
/usr/bin/lz4 -q -z %p /pgdata/10/archive_wals/%f.lz4
(1 row)

```

13.2.2 备份策略

PostgreSQL 提供了 `pg_start_backup`、`pg_stop_backup` 两个 API，
 PostgreSQL 9.1 版本开始，
`pg_basebackup` 提供了基于流复制的备份策略，
`pg_basebackup` 支持 `tar` 归档，
 使用 `pg_start_backup`、`pg_stop_backup` 两个 API，
 使用 `rsync`、`scp` 等工具进行备份，
 API 提供了基于 PIR 的备份策略，
 使用 `pg_start_backup`、`pg_stop_backup` 两个 API。

1. 基于 API 的备份策略

使用 API 进行备份，
`pg_start_backup` 用于启动备份，
`pg_stop_backup` 用于停止备份。

1 使用 `pg_start_backup` 启动备份

```
pg_start_backup[...]
```

```
1[...WAL[...]
```

```
[...WAL[...]
```

NOTICE: WAL archiving is not enabled; you must ensure that all required WAL segments are copied through other means to complete the backup

```
[...WAL[...]  
WAL[...WAL[...]  
pg_start_backup[...WAL[...]  
WAL[...13.2.1[...
```

```
2[...]
```

```
[...full_page_writes[...  
off[...full_page_writes[...  
off[...full_page_writes[...on[...  
[...]
```

```
3[...]
```


4 backup_label
backup_label

·START WAL LOCATION 25/2B002118
file 00000001000000250000002B

·CHECKPOINT LOCATION
LSN

·BACKUP METHOD
pg_start_backup pg_basebackup
BACK up METHOD streamed

·BACKUP FROM master
standby

·START TIME pg_start_backup

·LABEL pg_start_backup

pg_start_backup

```
pg_start_backup(label text [, fast boolean [, exclusive  
boolean ]])
```

label

```
pg_start_backup() CHECKPOINT()
fast() false() exclusive()
pg_start_backup()
()
() false() pg_is_in_backup()
()
() exclusive()
TRUE()
```

2 数据备份和恢复

```
rsync tar cp scp()
pg_wal pg_replslot()
postmaster.opts() postmaster.pid()
()
```

3 pg_stop_backup()

```
pg_stop_backup()
```

- pg_start_backup() full-page-writes() pg-stop-backup()

- XLOG()

- WAL()

· backup_label
pg_stop_backup

· backup_label backup_label
pg_start_backup

1 pg_start_backup

```
mydb=# SELECT pg_start_backup('base',false,false);
       pg_start_backup
-----
      0/4000028
(1 row)
```

2

```
[postgres@pghost1 ~]$ cd /pgdata/10/backups
[postgres@pghost1 /pgdata/10/backup]$ tar -cvf base.tar.gz
/pgdata/10/data --exclude=/pgdata/10/data/postmaster.pid --
exclude=/pgdata/10/data/postmaster.opts --
exclude=/pgdata/10/data/log/*
tar: Removing leading '/' from member names
/pgdata/10/data/
...
...
...
/pgdata/10/data/pg_wal/0000000100000000000000002.00000028.bac
kup
[postgres@pghost1 /pgdata/10/backup]$
[postgres@pghost1 /pgdata/10/backup]$ ll -h
```

```
total 80M
-rw-r--r-- 1 postgres postgres 80M Feb 13 00:35 base.tar.gz
```

3. pg_stop_backup

```
mydb=# SELECT pg_stop_backup(false);
NOTICE:  pg_stop_backup complete, all required WAL segments
have been archived
      pg_stop_backup
-----
(0/4000168,"START WAL LOCATION: 0/4000028 (file
00000001000000000000000004) +
CHECKPOINT LOCATION: 0/4000098
+
BACKUP METHOD: streamed
+
BACKUP FROM: master
+
START TIME: 2018-02-13 00:34:24 CST
+
LABEL: base
+
",")
(1 row)
```

pg_stop_backup

2. pg_basebackup

```
pg_basebackup -D /var/lib/postgresql/12/base -P
pg_basebackup -D /var/lib/postgresql/12/base -P
pg_basebackup -D /var/lib/postgresql/12/base -P
```

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_basebackup -Ft -
Pv -Xf -z -Z5 -p 1922 -D /pgdata/10/backups/
pg_basebackup: initiating base backup, waiting for
checkpoint to complete
pg_basebackup: checkpoint completed
pg_basebackup: write-ahead log start point: 33/A8000028 on
timeline 1
52782/52782 kB (100%), 1/1 tablespace
pg_basebackup: write-ahead log end point: 33/A8000130
pg_basebackup: base backup completed
```

base backup completed

```
[postgres@pghost1 ~]$ ll -h /pgdata/10/backups/
total 3.8M
-rw-r--r-- 1 postgres postgres 3.8M Feb 11 01:26 base.tar.gz
```

13.3 WAL

WAL

```
CREATE TABLE tbl
(
  id SERIAL PRIMARY KEY,
  ival INT NOT NULL DEFAULT 0,
  description TEXT,
  created_time TIMESTAMPTZ NOT NULL DEFAULT now()
);
```

```
mydb=# INSERT INTO tbl (ival) VALUES (1);
INSERT 0 1
mydb=# SELECT id,ival,description,created_time FROM tbl;
 id | ival | description |          created_time
-----+-----+-----+-----
   1 |    1 |              | 2018-02-13 01:26:36.767887+08
(1 row)
```

WAL 16MB

```
mydb=# SELECT pg_switch_wal();
pg_switch_wal
-----
0/3000350
(1 row)
```

archive_timeout timeout
WAL

13.3.1

ival 2
2

```
mydb=# INSERT INTO tbl (ival) VALUES (2);
INSERT 0 1
mydb=# SELECT id,ival,description,created_time FROM tbl;
 id | ival | description |      created_time
-----+-----+-----+-----
  1 |    1 |              | 2018-02-13 01:26:36.767887+08
  2 |    2 |              | 2018-02-13 01:47:42.280831+08
(2 rows)
```

created_time 01 47
42.280831+08

1
2

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_ctl -D
/pgdata/10/data/ -m stop
waiting for server to shut down.... done
server stopped
[postgres@pghost1 ~]$ rm -rf /pgdata/10/data
```

2 pg_basebackup

```
[postgres@pghost1 ~]$ mkdir -p /pgdata/10/data
[postgres@pghost1 ~]$ chmod 0700 /pgdata/10/data
[postgres@pghost1 ~]$ tar -xvf /pgdata/10/backups/base.tar.gz
-C /pgdata/10/data/
```

3 recovery.conf

share/recovery.conf
recovery.conf
0600

```
[postgres@pghost1 ~]$ cp /usr/pgsql-
10/share/recovery.conf.sample /pgdata/10/data/recovery.conf
[postgres@pghost1 ~]$ chmod 0600 /pgdata/10/data/recovery.conf
[postgres@pghost1 ~]$ ll /pgdata/10/data/recovery.conf
-rw----- 1 postgres postgres 5762 Feb 13 01:07
/pgdata/10/data/recovery.conf
```

recovery.conf

```
[postgres@pghost1 ~]$ vim /pgdata/10/data/recovery.conf
```

```
restore_command = 'ls -l /pgdata/10/archive_wals/%f.lz4 %p'

```

```
restore_command = '/usr/bin/lz4 -d
/pgdata/10/archive_wals/%f.lz4 %p'

```

```
recovery_target_timeline = 'latest'
recovery.conf

```

```
restore_command = '/usr/bin/lz4 -d
/pgdata/10/archive_wals/%f.lz4 %p'
recovery_target_timeline = 'latest'

```

4

```
recovery.conf

```

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_ctl -D
/pgdata/10/data start
[postgres@pghost1 ~]$
[postgres@pghost1 ~]$ tailf /pgdata/10/data/log/postgresql-
Tue.csv
...
...
...
2018-02-13 01:32:06.745 CST,,,41781,,5a81cf96.a335,2,,2018-02-
13 01:32:06 CST,,0,LOG,00000,"starting archive
recovery",,,,,,,,,,""
...
...
```

```

...
2018-02-13 01:32:06.908 CST,,,41781,,5a81cf96.a335,9,,2018-02-
13 01:32:06 CST,1/0,0,LOG,00000,"redo done at
0/3000060",,,,,,,,,,""
2018-02-13 01:32:06.946 CST,,,41781,,5a81cf96.a335,10,,2018-
02-13 01:32:06 CST,1/0,0,LOG,00000,"restored log file
""0000000100000000000000003"" from archive",,,,,,,,,,""
2018-02-13 01:32:06.961 CST,,,41781,,5a81cf96.a335,11,,2018-
02-13 01:32:06 CST,1/0,0,LOG,00000,"selected new timeline ID:
2",,,,,,,,,,""
2018-02-13 01:32:07.007 CST,,,41781,,5a81cf96.a335,12,,2018-
02-13 01:32:06 CST,1/0,0,LOG,00000,"archive recovery
complete",,,,,,,,,,""
2018-02-13 01:32:07.413 CST,,,41779,,5a81cf95.a333,3,,2018-02-
13 01:32:05 CST,,0,LOG,00000,"database system is ready to
accept connections",,,,,,,,,,""
...
...
...

```

recovery.conf PostgreSQL
 recovery.done

5

```

mydb=# SELECT id,ival,description,created_time FROM tbl;
   id | ival | description |          created_time
-----+-----+-----+-----
    1 |    1 |             | 2018-02-13 01:26:36.767887+08
    2 |    2 |             | 2018-02-13 01:47:42.280831+08
(2 rows)

```

13.3.2

mysql> SELECT * FROM tbl ORDER BY created_time DESC;

mysql> SELECT current_timestamp;

```
mydb=# SELECT id,ival,description,created_time FROM tbl ORDER
BY created_time DESC;
```

id	ival	description	created_time
3	3		2018-02-28 12:20:05.003997+08
2	2		2018-02-28 12:14:05.248928+08
1	1		2018-02-28 12:13:59.472933+08

(3 rows)

```
mydb=# SELECT current_timestamp;
current_timestamp
```

current_timestamp
2018-02-28 15:12:12.185701+08

(1 row)

mysql> DELETE FROM tbl WHERE

```
mydb=# DELETE FROM tbl;
DELETE 3
mydb=# SELECT id,ival,description,created_time FROM tbl ORDER
BY created_time DESC;
```

id	ival	description	created_time
----	------	-------------	--------------

(0 rows)

mysql> ALTER TABLE tbl ADD COLUMN current_timestamp

varchar(255) USING utf8;

```
restore_command = '/usr/bin/lz4 -d
/pgdata/10/archive_wals/%f.lz4 %p'
recovery_target_time = '2018-02-28 15:12:12.185701+08'
```

13.3.3 〇〇〇〇〇〇〇〇

pg_catalog.pg_create_restore_point
pg_catalog.pg_create_restore_point
pg_catalog.pg_create_restore_point

pg_catalog.pg_create_restore_point
pg_catalog.pg_create_restore_point

```
mydb=# \df pg_create_restore_point
```

Schema	Name	List of functions	Result data type
Argument data types	Type		
pg_catalog	pg_create_restore_point	pg_lsn	
text	normal		

(1 row)

pg_catalog.pg_create_restore_point
pg_catalog.pg_create_restore_point

```
mydb=# INSERT INTO tbl (ival) VALUES (4);
```

```
INSERT 0 1
```

```
mydb=# SELECT id,ival,description,created_time FROM tbl ORDER  
BY created_time DESC;
```

id	ival	description	created_time
4	4		2018-02-28 15:21:21.960729+08
3	3		2018-02-28 12:20:05.003997+08
2	2		2018-02-28 12:14:05.248928+08
1	1		2018-02-28 12:13:59.472933+08

(4 rows)

pg_catalog.pg_create_restore_point

```
mydb=# SELECT pg_create_restore_point('restore_point');
         pg_create_restore_point
-----
0/C0000C8
(1 row)
```

DROP

```
mydb=# DELETE FROM tbl WHERE id = 4;
DELETE 1
mydb=# ALTER TABLE tbl DROP COLUMN description;
ALTER TABLE
mydb=# SELECT * FROM tbl ORDER BY created_time DESC;
   id | ival |          created_time
-----+-----+-----
    3 |    3 | 2018-02-28 12:20:05.003997+08
    2 |    2 | 2018-02-28 12:14:05.248928+08
    1 |    1 | 2018-02-28 12:13:59.472933+08
(3 rows)
```

“restore_point”

```
[postgres@pghost1 /pgdata/10/data]$ /usr/pgsql-10/bin/pg_ctl -
D /pgdata/10/data -mi stop
waiting for server to shut down.... done
server stopped
[postgres@pghost1 /pgdata/10/data]$
[postgres@pghost1 /pgdata/10/data]$ rm -rf *
[postgres@pghost1 /pgdata/10/data]$ tar -xf
/pgdata/10/backups/base.tar.gz -C /pgdata/10/data/
[postgres@pghost1 /pgdata/10/data]$ cp /usr/pgsql-
10/share/recovery.conf.sample /pgdata/10/data/recovery.conf
[postgres@pghost1 /pgdata/10/data]$ chmod 0600
/pgdata/10/data/recovery.conf
[postgres@pghost1 /pgdata/10/data]$ ll
/pgdata/10/data/recovery.conf
```

```
-rw----- 1 postgres postgres 5762 Feb 13 01:07
/pgdata/10/data/recovery.conf
```

recovery.conf

```
[postgres@pghost1 ~]$ vim /pgdata/10/data/recovery.conf
```

restore_command
recovery_target_name

```
restore_command = '/usr/bin/lz4 -d
/pgdata/10/archive_wals/%f.lz4 %p'
recovery_target_name = 'restore_point'
```

```
...
...
...
2018-02-28 16:28:33.396 CST,,,32457,,5a966831.7ec9,2,,2018-02-
28 16:28:33 CST,,0,LOG,00000,"starting point-in-time recovery
to ""restore_point""",,,,,,""
...
...
...
2018-02-28 16:28:33.450 CST,,,32455,,5a966830.7ec7,2,,2018-02-
28 16:28:32 CST,,0,LOG,00000,"database system is ready to
accept read only connections",,,,,,""
...
...
...
2018-02-28 16:28:33.985 CST,,,32455,,5a966830.7ec7,3,,2018-02-
28 16:28:32 CST,,0,LOG,00000,"database system is ready to
accept connections",,,,,,""
```

...

restore_point

```
mydb=# SELECT * FROM tbl ORDER BY created_time DESC;
 id | ival | description |          created_time
-----+-----+-----+-----
   4 |    4 |              | 2018-02-28 15:21:21.960729+08
   3 |    3 |              | 2018-02-28 12:20:05.003997+08
   2 |    2 |              | 2018-02-28 12:14:05.248928+08
   1 |    1 |              | 2018-02-28 12:13:59.472933+08
(4 rows)
```

restore_point

13.3.4

PostgreSQL

restore_point

```
mydb=# SELECT * FROM tbl ORDER BY created_time DESC;
 id | ival | description |          created_time
-----+-----+-----+-----
   4 |    4 |              | 2018-02-28 15:21:21.960729+08
   3 |    3 |              | 2018-02-28 12:20:05.003997+08
   2 |    2 |              | 2018-02-28 12:14:05.248928+08
   1 |    1 |              | 2018-02-28 12:13:59.472933+08
(4 rows)
```


pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1 -Ft -v -j1 -Z lz4 -T 1 -t 1

```
mydb=# BEGIN;
BEGIN
mydb=# SELECT txid_current();
      txid_current
-----
          561
(1 row)
mydb=# DELETE FROM tbl WHERE id > 1;
DELETE 3
mydb=# END;
COMMIT
```

pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1

```
mydb=# SELECT * FROM tbl ORDER BY created_time DESC;
   id | ival | description |      created_time
-----+-----+-----+-----
    1 |    1 |              | 2018-02-28 12:13:59.472933+08
(1 row)
```

pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1
pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1
pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1

```
restore_command = '/usr/bin/lz4 -d  
/pgdata/10/archive_wals/%f.lz4 %p'  
recovery_target_xid = 561
```

pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1
pg_restore -c -d -Ft -v -j1 -Z lz4 -T 1 -t 1

□□□□□□□□

```
...
...
...
2018-02-28 18:06:55.830 CST,,,7036,,5a967f3f.1b7c,2,,2018-02-
28 18:06:55 CST,,0,LOG,00000,"starting point-in-time recovery
to XID 561",,,,,,,,,,""
    2018-02-28 18:06:56.003 CST,,,7036,,5a967f3f.1b7c,9,,2018-
02-28 18:06:55 CST,1/0,0,LOG,00000,"selected new timeline ID:
7",,,,,,,,,,""
...
...
...
2018-02-28 18:06:55.883 CST,,,7033,,5a967f3e.1b79,2,,2018-02-
28 18:06:54 CST,,0,LOG,00000,"database system is ready to
accept read only connections",,,,,,,,,,""
...
...
...
```

□□□□□□□□□□□□□□□□□□

```
mydb=# SELECT * FROM tbl ORDER BY created_time DESC;
   id | ival | description |          created_time
-----+-----+-----+-----
    4 |    4 |              | 2018-02-28 15:21:21.960729+08
    3 |    3 |              | 2018-02-28 12:20:05.003997+08
    2 |    2 |              | 2018-02-28 12:14:05.248928+08
    1 |    1 |              | 2018-02-28 12:13:59.472933+08
(4 rows)
```

□□□□□□□□□□□□□□□□□□

13.3.5 □□□□□□□□

pg13.3.1 recovery_target_timeline timeline

initdb TimeLineID 1 TimeLineID

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/initdb -D /pgdata/10/data
...
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_controldata /pgdata/10/data | grep TimeLineID
Latest checkpoint's TimeLineID:      1
Latest checkpoint's PrevTimeLineID:  1
```

TimeLineID 1

```
mydb=# INSERT INTO tbl (ival,description) VALUES (1,'');
INSERT 0 1
mydb=# INSERT INTO tbl (ival,description) VALUES (2,'');
INSERT 0 1
mydb=# SELECT id,ival,description,created_time FROM tbl;
 id | ival | description |      created_time
-----+-----+-----+-----
   1 |    1 |          | 2018-02-28 21:23:41.881687+08
   2 |    2 |          | 2018-02-28 21:24:25.07764+08
(2 rows)
mydb=# SELECT pg_switch_wal();
pg_switch_wal
```

0/167D4D8
(1 row)

pg_wal

```
[postgres@pghost1 ~]$ ll -h /pgdata/10/data/pg_wal/  
total 33M  
-rw----- 1 postgres postgres 16M Feb 28 21:29  
00000001000000000000000001  
-rw----- 1 postgres postgres 16M Feb 28 21:29  
00000001000000000000000002  
drwx----- 2 postgres postgres 4.0K Feb 28 21:29  
archive_status  
[postgres@pghost1 ~]$ ll -h  
/pgdata/10/data/pg_wal/archive_status/  
total 0  
-rw----- 1 postgres postgres 0 Feb 28 21:29  
00000001000000000000000001.done  
[postgres@pghost1 ~]$ ll -h /pgdata/10/archive_wals/  
total 2.3M  
-rw----- 1 postgres postgres 2.3M Feb 28 21:29  
00000001000000000000000001.lz4
```

pg_switch_wal WAL
pg_wal
archive_status

TimeLineID

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_controldata  
/pgdata/10/data | grep TimeLineID  
Latest checkpoint's TimeLineID: 2  
Latest checkpoint's PrevTimeLineID: 1
```

pg_wal

```
[postgres@pghost1 ~]$ ll -h /pgdata/10/data/pg_wal/
total 33M
-rw----- 1 postgres postgres 16M Feb 28 22:19
00000001000000000000000003
-rw----- 1 postgres postgres 16M Feb 28 22:21
00000002000000000000000004
-rw----- 1 postgres postgres 41 Feb 28 22:19
00000002.history
drwx----- 2 postgres postgres 4.0K Feb 28 22:19
archive_status
[postgres@pghost1 ~]$ ll -h
/pgdata/10/data/pg_wal/archive_status/
total 0
-rw----- 1 postgres postgres 0 Feb 28 21:33
00000001000000000000000003.done
-rw----- 1 postgres postgres 0 Feb 28 22:19
00000002.history.done
```

pg_wal WAL

“.history”

TimeLineID TimeLineID 1 4

“TimeLineID.history”

TimeLineID 2

00000002.history

```
[postgres@pghost1 ~]$ ll -h /pgdata/10/archive_wals/
total 2.5M
-rw----- 1 postgres postgres 2.3M Feb 28 21:29
```

```
0000000100000000000000001.lz4
-rw----- 1 postgres postgres 78K Feb 28 21:33
0000000100000000000000002.lz4
-rw----- 1 postgres postgres 216 Feb 28 21:33
0000000100000000000000003.00000028.backup.lz4
-rw----- 1 postgres postgres 78K Feb 28 21:33
0000000100000000000000003.lz4
-rw----- 1 postgres postgres 60 Feb 28 22:19
00000002.history.lz4
```

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

```
mydb=# INSERT INTO tbl (ival,description) VALUES (3,'XXXXXX
X');
INSERT 0 1
mydb=# INSERT INTO tbl (ival,description) VALUES (4,'XXXXX');
INSERT 0 1
mydb=# SELECT id,ival,description,created_time FROM tbl;
 id | ival | description |          created_time
-----+-----+-----+-----
  1 |    1 | XXXXXXXX | 2018-02-28 21:23:41.881687+08
  2 |    2 | XXXXXX | 2018-02-28 21:24:25.07764+08
  3 |    3 | XXXXXXXX | 2018-02-28 22:21:03.68966+08
  4 |    4 | XXXXX | 2018-02-28 23:01:58.253911+08
(4 rows)
```

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX"4|4|XXXX

X|2018-02-2823X01X58.253911+08"XXXX

XXXXXXXXXXXXrecovery_target_timelineXXXX2XXXX

XXXXXXXXXXXXrecovery.confXXXXXX

```
restore_command = '/usr/bin/lz4 -d
/pgdata/10/archive_wals/%f.lz4 %p'
recovery_target_timeline = 2
recovery_target_time = '2018-02-28 23:00:00'
```

□□□□□□□□□□□□□□

```
mydb=# SELECT id,ival,description,created_time FROM tbl;
 id | ival | description | created_time
-----+-----+-----+-----
-
   1 |    1 | □□□□□□□ | 2018-02-28 21:23:41.881687+08
   2 |    2 | □□□□□□ | 2018-02-28 21:24:25.07764+08
   3 |    3 | □□□□□□□ | 2018-02-28 22:21:03.68966+08
(3 rows)
```

□□□□□□□□□□□□2□□□□□□□□

13.4 SQL

SQL is a standard language for interacting with databases. It is used to create, modify, and retrieve data from a database. SQL is a declarative language, meaning that you specify what you want, not how to get it.

PostgreSQL is a powerful, open source object-relational database system. It is known for its reliability, security, and performance. PostgreSQL supports SQL, and it provides a rich set of features for managing data.

13.4.1 SQL

SQL is a standard language for interacting with databases. It is used to create, modify, and retrieve data from a database. SQL is a declarative language, meaning that you specify what you want, not how to get it. PostgreSQL is a powerful, open source object-relational database system. It is known for its reliability, security, and performance. PostgreSQL supports SQL, and it provides a rich set of features for managing data. The `pg_dump` and `pg_dumpall` utilities are used to create backups of PostgreSQL databases. `pg_dump` is used to dump a single database, and `pg_dumpall` is used to dump all databases in a cluster. Both utilities support various output formats, including plain text, custom, and directory.

1. `pg_dump` and `pg_dumpall`

`pg_dump`

`pg_dump` dumps a database as a text file or to other formats.
Usage:

-c, --clean

Drop database

-C, --create

Create database

-n, --schema=SCHEMA
-N, --exclude-schema=SCHEMA

Schema Schema

-s, --schema-only
-t, --table=TABLE
-T, --exclude-table=TABLE

Schema

--inserts	dump data as INSERT commands,
rather than COPY	
--column-inserts	dump data as INSERT commands
with column names	

pg_dump SQL PostgreSQL
COPY SQL

INSERT INTO INSERT
PostgreSQL
--insert--column-inserts
INSERT

SQL
INSERT

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_dump -Fp -a --
insert --column-inserts -t tbl -p 1921 mydb > dump.sql
[postgres@pghost1 ~]$ cat dump.sql
--
-- PostgreSQL database dump
--
-- Dumped from database version 10.2
-- Dumped by pg_dump version 10.2
SET statement_timeout = 0;
...
...
SET search_path = public, pg_catalog;
--
-- Data for Name: tbl; Type: TABLE DATA; Schema: public;
Owner: postgres
--
INSERT INTO tbl (id, ival, description, created_time) VALUES
(1, 1, ' ', '2018-02-28 21:23:41.881687+08');
INSERT INTO tbl (id, ival, description, created_time) VALUES
(2, 2, ' ', '2018-02-28 21:24:25.07764+08');
INSERT INTO tbl (id, ival, description, created_time) VALUES
(3, 3, ' ', '2018-02-28 22:21:03.68966+08');
--
-- Name: tbl_id_seq; Type: SEQUENCE SET; Schema: public;
Owner: postgres
--
SELECT pg_catalog.setval('tbl_id_seq', 36, true);
--
```

```
-- PostgreSQL database dump complete
--
```

pg_dumpall pg_dump

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_dumpall -r -p
1921
```

```
--
-- PostgreSQL database cluster dump
--
SET default_transaction_read_only = off;
SET client_encoding = 'UTF8';
SET standard_conforming_strings = on;
--
-- Roles
--
CREATE ROLE postgres;
ALTER ROLE postgres WITH SUPERUSER INHERIT CREATEROLE
CREATEDB LOGIN REPLICATION BYPASSRLS;
--
-- PostgreSQL database cluster dump complete
--
```

2.SQL

psql SQL

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/psql -p 1922 mydb <
dump.sql
```

13.5

```

pg_rman[pg_backrest]

```

14

Oracle DBA Oracle RAC PostgreSQL RAC PostgreSQL 12

```

PostgreSQL
WAL
VIP
virtual ip
IP
VIP

```

A diagram showing a 6x24 grid of squares. The top row has the word "WALL" in large black letters, with the first 18 squares being white and the last 6 squares being black. The remaining 5 rows are entirely black.

数据库连接池
数据库连接池

数据库连接池PostgreSQL数据库连接池
Pgpool-II+数据库连接池Keepalived+数据库连接池

14.1 Pgpool-II+□□□□□□□□□□

□□□□□Pgpool-II□□□□□□□□□□□□□Pgpool-
 II□□□□□□□□□□□□□□□□□□□□□□□□□□□□
 □□Pgpool-II□□□□□□□□□□□□□□□Pgpool-II□□□
 pgpool□

```

      ·pgpool
pgpool

```

```
· SELECT
pgpool
```

```
·[ ] [ ] [ ] [ ] [ ] [ ] [ ] [ ] [ ] [ ] pgpool [ ] [ ] [ ] [ ] [ ] [ ]  
[ ] [ ] [ ] [ ] [ ] [ ]
```

```
·pgpool PostgreSQL
pgpool
PostgreSQL
```

```

pgpool
pgpool
pgpool

```

pgpoolをインストールする
pgpoolをインストールする

・PostgreSQLをインストールする
PostgreSQLをインストールする
pgpoolをインストールする
\$prefix/etc/pgpool.conf.sample-stream
をインストールする

・Slonyをインストールする
pgpoolをインストールする
Slonyをインストールする
\$prefix/etc/pgpool.conf.sample-master-slave
をインストールする

・pgpoolをインストールする
pgpoolをインストールする
pgpoolをインストールする
\$prefix/etc/pgpool.conf.sample-replication
をインストールする

・pgpoolをインストールする
pgpoolをインストールする
pgpoolをインストールする
\$prefix/etc/pgpool.conf.sample
をインストールする

Pgpoolをインストールする
Pgpoolをインストールする

搭建基于PostgreSQL12.1的14-1

14-1 pgpool+搭建

主 机 名	组 件	IP 地址	端 口	版 本	VIP
pghost4	主库	192.168.26.57	1921	PostgreSQL10	192.168.26.72
	pgpool 主	192.168.26.57	9999	Pgpool-II 3.6.6	
pghost5	备库	192.168.26.58	1921	PostgreSQL10	
	pgpool 备	192.168.26.58	9999	Pgpool-II 3.6.6	

14.1.1 pgpool搭建

pgpool14-1

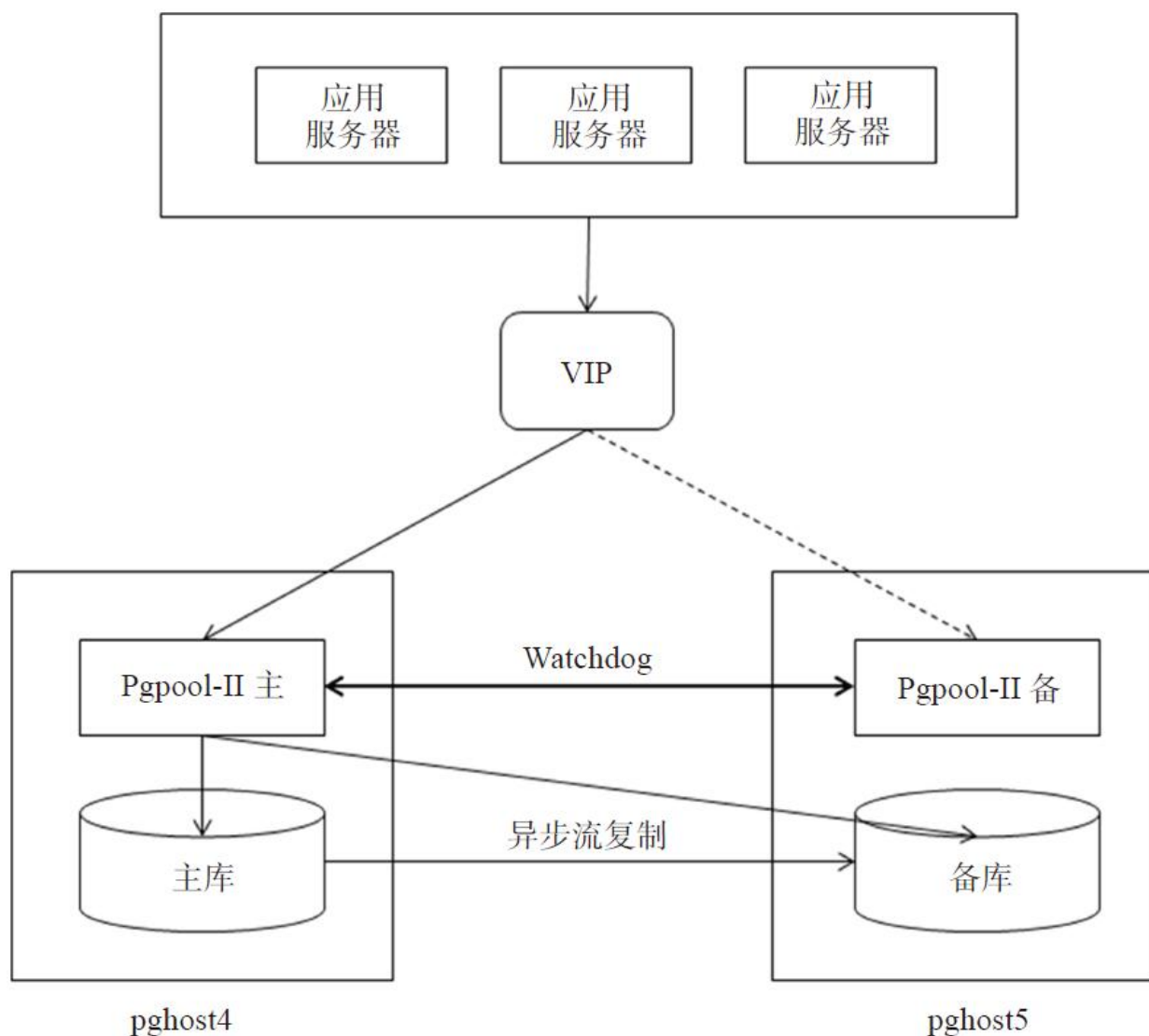


图14-1 Pgpool-II架构图

pghost4、pghost5、Pgpool-II、Pgpool-II、pghost4、
 Pgpool-II、pghost5、Pgpool-II、Pgpool-II、
 pghost4、Pgpool-II、pghost5、Pgpool-II、
 Watchdog、Pgpool-II、Pgpool-II、Pgpool-II、
 Pgpool-II、Pgpool-II

watchdog pgpool watchdog
pgpool pgpool
watchdog

· pgpool PostgreSQL
pgpool

· pgpool

· watchdog pgpool
pgpool pgpool
pgpool

· pgpool watchdog
watchdog

· pgpool watchdog VIP
pgpool

14.1.2 pgpool

pgpool Pgpool-II 3.6.6
RPM
[http://www.pgpool.net/download.php?
f=pgpool-II-3.6.6.tar.gz](http://www.pgpool.net/download.php?f=pgpool-II-3.6.6.tar.gz)

pgghost4 pgghost5 pgpool-II-3.6.6.tar.gz

```
# tar xvf pgpool-II-3.6.6.tar.gz
```

```
# ./configure --prefix=/opt/pgpool --with-pgsql=/opt/pgsql
# make
# make install
```

/opt/pgpool

```
[root@pgghost4 etc]# ll /opt/pgpool/
total 20
drwxr-xr-x 2 root root 4096 Oct 4 14:39 bin
drwxr-xr-x 2 root root 4096 Oct 4 19:13 etc
drwxr-xr-x 2 root root 4096 Oct 4 14:39 include
drwxr-xr-x 2 root root 4096 Oct 4 14:39 lib
drwxr-xr-x 3 root root 4096 Oct 4 14:39 share
```

/opt/pgpool/bin pgpool
pgpool pg_md5 pcp_attach_node
/opt/pgpool/etc
pcp.conf.sample pgpool.conf.sample-stream

root@pgpool:~# cd /opt/pgpool/pgpool
root@pgpool:~#

1. pghost4 pghost5 /etc/hosts

/etc/hosts

```
192.168.26.57 pghost4
192.168.26.58 pghost5
```

2. pghost4 pghost5

root@pgpool:~# cd /opt/pgpool/pgpool
pgpool/pgpool
failover_command ssh
pgpool

pghost4 pghost5 postgres
pgpool/pgpool pghost4

```
# su - postgres
$ ssh-keygen
$ ssh-copy-id postgres@pghost5
```

pghost4 postgres
ssh pghost5

```
ssh postgres@pghost5
```

```
cd /opt/pgghost5
```

```
cd /opt/pgpool
```

3.pgghost4pgghost5pool_hba.conf

```
PostgreSQL
$PGDATA/pg_hba.conf
pgpool PostgreSQL
pgpool pgpool
pgpool pgpool
```

```
pool_hba.conf.sample
pool_hba.conf
```

```
# cd /opt/pgpool/etc
# cp pool_hba.conf.sample pool_hba.conf
```

```
pool_hba.conf
PostgreSQL pg_hba.conf
pool_hba.conf
```

host	replication	repuser	192.168.26.57/32
md5			
host	replication	repuser	192.168.26.58/32
md5			

host	replication	repuser	192.168.26.72/32
md5			
host	all	all	0.0.0.0/0 md5

4.pgghost4pgghost5pool_passwd

```
pgpool MD5 pgpool
pool_passwd pool_passwd
```

```
username: encrypted_passwd
```

```
pgpool pg_md5
```

```
# pg_md5 -u postgres -m postgres123
```

```
pg_md5 MD5 MD5
pool_passwd
```

```
# cat pool_passwd
postgres:md5163311300b0732b814a34aabfdfffe62
```

```
postgres
pool_passwd
```

```
postgres=# SELECT rolpassword FROM pg_authid WHERE
rolname='postgres';
           rolpassword
-----
md5163311300b0732b814a34aabfdfffe62
(1 row)
```

postgres角色密码MD5
pool_passwd

5.pgghost4pgghost5pgpool.conf

/opt/pgpool/etc

```
# cd /opt/pgpool/etc
# cp pgpool.conf.sample-stream pgpool.conf
```

pgpool.conf

pgpool

```
listen_addresses = '*'
port = 9999
```

·listen_addressespgpool*
pgpool

·port=9999pgpool9999

pgpool PostgreSQL
pghost4 pgghost5 0
1

```
backend_hostname0 = 'pghost4'
backend_port0 = 1921
backend_data_directory0 = '/data1/pg10/pg_root'
backend_flag0 = 'ALLOW_TO_FAILOVER'
backend_hostname1 = 'pghost5'
backend_port1 = 1921
backend_data_directory1 = '/data1/pg10/pg_root'
backend_flag1 = 'ALLOW_TO_FAILOVER'
```

·backend_hostname0 PostgreSQL
0 IP

·backend_port0 PostgreSQL 0

·backend_data_directory0
PostgreSQL 0

·backend_flag0
ALLOW_TO_FAILOVER

·backend_hostname1 PostgreSQL
1 IP

·backend_port1 PostgreSQL 1

·backend_data_directory1
PostgreSQL1

·backend_flag1
ALLOW_TO_FAILOVER

pgpool

```
enable_pool_hba = on  
pool_passwd = 'pool_passwd'
```

·enable_pool_hba
pool_hba.conf

·pool_passwd
pool_passwd

pgpoolpid

```
log_destination = 'syslog'  
pid_file_name = '/opt/pgpool/pgpool.pid'
```

·log_destinationpgpool
stderrsyslog
syslog/var/log/messagepgpool

·pid_file_name pgpool PID

pgpool

```
load_balance_mode = off
```

```
load_balance_mode pgpool
SELECT pgpool

```

pgpool

```
# pgpool
master_slave_mode = on
master_slave_sub_mode = 'stream'
sr_check_period = 10
sr_check_user = 'repuser'
sr_check_password = 're12a345'
sr_check_database = 'postgres'
delay_threshold = 10000000
```

·master_slave_mode
off on

·master_slave_sub_mode
slony stream slony slony stream
PostgreSQL stream

·sr_check_period健康チェックの周期10
分

·sr_check_user健康チェックを行うユーザ

·sr_check_database健康チェックを行うデータベース

·delay_threshold健康チェックの遅延の閾値WAL
ログpgpoolがSELECTする際の遅延

pgpoolはPostgreSQLのクライアント
ライブラリを介してデータベースと通信する
ため、データベースの健康状態を確認する

```
health_check_period = 5
health_check_timeout = 20
health_check_user = 'repuser'
health_check_password = 're12a345'
health_check_database = 'postgres'
health_check_max_retries = 3
health_check_retry_delay = 3
```

健康チェックの失敗時の処理

```
failover_command = '/opt/pgpool/failover_stream.sh %d %P %H  
%R'
```

```
failover_command pgpool
pgpool
pgpool %d ID %P
ID %H %R
failover_stream.sh
```

watchdog

```
use_watchdog = on
wd_hostname = 'pgghost4'
wd_port = 9000
wd_priority = 1
```

```
·use_watchdog watchdog off
·wd_hostname watchdog IP
pgpool
·wd_port watchdog 9000
·wd_priority watchdog pgpool
watchdog pgpool
pgpool
```

VIP

```
delegate_IP = '192.168.26.72'
if_cmd_path = '/sbin'
```

```
if_up_cmd = 'ip addr add $_IP_$/24 dev bond0 label bond0:1'
if_down_cmd = 'ip addr del $_IP_$/24 dev bond0'
```

```
·delegate_IP pgpool VIP pgpool
pgpool VIP pgpool pgpool pgpool pgpool pgpool
pgpool pgpool VIP
```

```
·if_cmd_path pgpool VIP pgpool
```

```
·if_up_cmd pgpool VIP pgpool ip addr add
pgpool VIP pgpool4 pgpool5 pgpool pgpool
bond0 pgpool VIP pgpool pgpool bond0 1 pgpool
pgpool IP
```

```
·if_down_cmd pgpool VIP pgpool ip addr
del
```

```
watchdog pgpool
```

```
heartbeat_destination0 = 'pgpool5' # pgpool IP
heartbeat_destination_port0 = 9694 # pgpool
heartbeat_device0 = 'bond0'
```

```
·heartbeat_destination0 pgpool
pgpool IP watchdog pgpool
heartbeat_destination0
```


·heartbeat_destination_port0
pgpool9694

·heartbeat_device0pgpool
watchdogbond0

watchdog

```
wd_life_point = 3
wd_lifecycle_query = 'SELECT 1'
wd_lifecycle_dbname = 'postgres'
wd_lifecycle_user = 'repuser'
wd_lifecycle_password = 're12a345'
```

·wd_life_pointpgpool
00

·wd_lifecycle_querypgpool
SQL

·wd_lifecycle_dbnamepgpool
0000

·wd_lifecycle_userpgpool
0000

·wd_lifecycle_passwordpgpool
00000000

pgpool 設定ファイル

```
other_pgpool_hostname0 = 'pghost5'    # pgpoolのIP  
other_pgpool_port0 = 9999             # pgpoolのポート  
other_wd_port0 = 9000                 # pgpoolのwatchdog  
#
```

・other_pgpool_hostname0 pgpool
のIP

・other_pgpool_port0 pgpool
のポート

・other_wd_port0 pgpool
watchdogのポート

pgpool.confのデフォルト設定は
pgpoolの
[http://www.pgpool.net/docs/latest/en/html/
runtime-config.html](http://www.pgpool.net/docs/latest/en/html/runtime-config.html) 参照



pgpoolのpgpool.conf
pgpool
pgpool

/etc/pgpool-II/failover_stream.sh
参照

```
#!/bin/bash
# Execute command by failover.
# special values: %d = node id
#                  %h = host name
#                  %p = port number
#                  %D = database cluster path
#                  %m = new master node id
#                  %M = old master node id
#                  %H = new master node host name
#                  %P = old primary node id
#                  %R = new master database cluster path
#                  %r = new master port number
#                  %% = '%' character

falling_node=$1          # %d
old_primary=$2           # %P
new_primary=$3           # %H
pgdata=$4               # %R
pghome=/opt/pgsql
log=/tmp/failover.log
date >> $log

# 
echo "falling_node=$falling_node" >> $log
echo "old_primary=$old_primary" >> $log
echo "new_primary=$new_primary" >> $log
echo "pgdata=$pgdata" >> $log

# root
if [ $falling_node = $old_primary ] && [ $UID -eq 0 ]; then
# $PGDATA/recovery.conf
    if [ -f $pgdata/recovery.conf ]; then
        su postgres -c "$pghome/bin/pg_ctl promote -D $pgdata"
        echo "Local promote" >> $log
    else
        su postgres -c "ssh -T postgres@$new_primary $pghome/bin/pg_ctl promote -D $pgdata"
        echo "Remote promote" >> $log
    fi
fi;
exit 0;
```

```
pgpool
pgpool
pgpool
$PGDATA
recovery.conf
pgpool
```

```
pgpool ip addr
IP root pgpool root
```

```
export PGP00L_HOME=/opt/pgpool
export PATH=$PGP00L_HOME/bin:$PATH:.
```

```
pgghost4 root pgpool
```

```
# pgpool
```

```

pgpoolpgpoolpgpool
pgpoolpgpoolpgpoolpgpool

```

```
·-a--hba-filepgpoolhba
$prefix/etc/pool hba.conf
```

```

    -f--config-filepgpool
$prefix/etc/pgpool.conf

```

```
.-F--pcp-filepgpoolpcp
$prefix/etc/pcp.conf
```

```

/opt/pgpool/etc
pgpool
```

```

pgpool.conf
reload
```

```
# pgpool reload
```

```
pgpool
```

```
# pgpool -m fast stop
```

```
pgpoolsmartfast
immediateSMART
pgpoolfastimmediate
pgpoolfast
```

```

pgghost4/var/log/messages
pgghost5
pgpoolpgghost5
/var/log/messagespgpool
```

☐VIP☐☐pgpool☐☐☐☐☐☐

```
[postgres@pghost4 ~]$ psql -h 192.168.26.72 -p 9999 postgres
postgres
postgres=# show pool_nodes;
   node_id | hostname | port | status | lb_weight | role
| select_cnt | load_balance_node | replication_delay
-----+-----+-----+-----+-----+-----+-----
---+-----+-----+-----+-----+-----+-----
      0    | pghost4  | 1921 | up     | -nan     | primary | 0
| true      | 0
      1    | pghost5  | 1921 | up     | -nan     | standby | 0
| false     | 0
(2 rows)
```

```

pgpool[0000000000]9999[0000000000]
pgpool[0000000000]status[0000000000]up[00]
pgpool[0000000000]down[00]pgpool[0000000000]
role[0000000000]primary[000000]standby[0000]
[00]node_id[00]0[0000000000]pghost4[00000000]
node_id[00]1[0000000000]pghost5[00000000]

```

[illegible]

14.1.3 PCP問題

`pgpool` `pgpool` `pgpool` `pgpool`

PCPはpgpoolをPostgreSQL
に接続するためのpcpのMD5パスワードを
pcp.confにpcp.confに設定する

```
# USERID:MD5PASSWD
```

USERIDのMD5PASSWORDを
設定する

pgpoolのPCPパスワードを
PCPのpgpoolのpg_md5のMD5
パスワード

```
# pg_md5 pgpool  
ba777e4c2f15c11ea8ac3be7e0440aa0
```

pcp.confに

```
# USERID:MD5PASSWD  
pgpool:ba777e4c2f15c11ea8ac3be7e0440aa0
```

PCPのpcp_node_info
pcp_watchdog_infoのpcp_attach_node
のpcp_node_infoを

```
$ pcp_node_info --verbose -h 192.168.26.72 -U pgpool 0
Hostname      : pghost4
Port          : 1921
Status        : 2
Weight        : nan
Status Name: up
```

0 Status

• 0 PCP

• 1

• 2

• 3

pcp_attach_node

14.1.4 pgpool

pgpool

pgpool pgpool

•

• 配置pgpool

pgpool 配置

pgghost4 pgpool pgghost5
pgpool VIP pgghost4

pcp_watchdog_info pgpool
watchdog watchdog pgpool
watchdog pgpool
watchdog pgpool

```
[postgres@pgghost4 ~]$ pcp_watchdog_info --verbose -h  
192.168.26.72 -U pgpool
```

```
Password:
```

```
Watchdog Cluster Information
```

```
Total Nodes      : 2  
Remote Nodes     : 1  
Quorum state     : QUORUM EXIST  
Alive Remote Nodes : 1  
VIP up on local node : YES  
Master Node Name  : pgghost4:9999 Linux pgghost4  
Master Host Name  : pgghost4
```

```
Watchdog Node Information
```

```
Node Name       : pgghost4:9999 Linux pgghost4  
Host Name       : pgghost4  
Delegate IP     : 192.168.26.72  
Pgpool port     : 9999  
Watchdog port   : 9000  
Node priority   : 1  
Status          : 4  
Status Name     : MASTER
```

```
Node Name       : pgghost5:9999 Linux pgghost5
```

Host Name : pghost5
Delegate IP : 192.168.26.72
Pgpool port : 9999
Watchdog port : 9000
Node priority : 1
Status : 7
Status Name : STANDBY

watchdog pghost4
watchdog pghost5 watchdog
pgpool watchdog

pghost4 IP

```
[postgres@pghost4 ~]$ ip a
...
12: bond0: <BROADCAST,MULTICAST,MASTER,UP,LOWER_UP> mtu 1500
    qdisc noqueue state UP
        link/ether a0:36:9f:9b:07:af brd ff:ff:ff:ff:ff:ff
        inet 192.168.26.57/24 brd 192.168.26.255 scope global
bond0
    inet 192.168.26.72/24 scope global secondary bond0:1
...
```

192.168.26.72 VIP pghost4
pghost4 pgpool

pghost4 pgpool

```
# pgpool -m fast stop
.done.
```

pgpool/var/log/messages

```
Oct 12 10:09:39 pghost4 pgpool[41818]: [13-1] 2017-10-12
10:09:39: pid 41818: LOG: received fast shutdown request
Oct 12 10:09:39 pghost4 pgpool[41818]: [14-1] 2017-10-12
10:09:39: pid 41818: LOG: shutdown request. closing listen
socket
Oct 12 10:09:39 pghost4 pgpool[26337]: [1-1] 2017-10-12
10:09:39: pid 26337: LOG: stop request sent to pgpool.
waiting for termination...
Oct 12 10:09:39 pghost4 pgpool: watchdog[41819]: [45-1]
2017-10-12 10:09:39: pid 41819: LOG: Watchdog is shutting
down
Oct 12 10:09:39 pghost4 pgpool: watchdog de-
escalation[26338]: [46-1] 2017-10-12 10:09:39: pid 26338:
LOG: watchdog: de-escalation started
Oct 12 10:09:39 pghost4 pgpool: watchdog de-
escalation[26338]: [47-1] 2017-10-12 10:09:39: pid 26338:
LOG: successfully released the delegate IP:"192.168.26.72"
Oct 12 10:09:39 pghost4 pgpool: watchdog de-
escalation[26338]: [47-2] 2017-10-12 10:09:39: pid 26338:
DETAIL: 'if_down_cmd' returned with success
Oct 12 10:09:41 pghost4 ntpd[6835]: Deleting interface #14
bond0:1, 192.168.26.72#123, interface stats: received=0,
sent=0, dropped=0, active_time=390 secs
```

pgpoolfast
watchdogpgpool.conf
if_down_cmdVIPDeleting
interfaceVIP

pghost4IPVIP

```
[postgres@pghost4 ~]$ ip a
12: bond0: <BROADCAST,MULTICAST,MASTER,UP,LOWER_UP> mtu 1500
qdisc noqueue state UP
    link/ether a0:36:9f:9b:07:af brd ff:ff:ff:ff:ff:ff
    inet 192.168.26.57/24 brd 192.168.26.255 scope global
    bond0
```

pgghost5 IP VIP pgghost5
bond0

pcp_watchdog_info pgpool
pgghost5 pgpool MASTER
pgghost4 pgpool SHUTDOWN

pgpool

```
[postgres@pghost4 ~]$ psql -h 192.168.26.72 -p 9999 postgres
postgres
Password for user postgres:
psql (10.0)
Type "help" for help.
```

pgpool

pgpool
PostgreSQL pgghost4
pgghost5 pgpool
pgpool
pgpool

pgpool 2.2.10 安装

pgghost4 和 pgpool 安装
pgghost5 和 pgpool 安装
VIP 和 pgghost4

pgpool 安装

```
[postgres@pgghost4 ~]$ psql -h 192.168.26.72 -p 9999 postgres
postgres=# show pool_nodes;
 node_id | hostname | port | status | lb_weight | role
 | select_cnt | load_balance_node | replication_delay
-----+-----+-----+-----+-----+-----
 0       | pgghost4 | 1921 | up     | -nan      | primary | 0
 | true      | 0
 1       | pgghost5 | 1921 | up     | -nan      | standby | 0
 | false     | 0
(2 rows)
```

pgghost4 安装

```
[postgres@pgghost4 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

pgpool 安装
pgpool.conf
failover_command
/opt/pgpool/failover_stream.sh
/tmp/failover.log

```
Thu Oct 12 11:01:50 CST 2017
falling_node=0
old_primary=0
new_primary=pghost5
pgdata=/data1/pg10/pg_root
Remote promote
```

falling_node PostgreSQL
old_primary new_primary
pgdata

pghost5

```
[postgres@pghost5 ~]$ pg_controldata | grep cluster
Database cluster state:          in production
```

pghost5 pghost5
pgdata

pgpool

```
[postgres@pghost4 ~]$ psql -h 192.168.26.72 -p 9999 postgres
postgres
postgres=# show pool_nodes;
 node_id | hostname | port | status | lb_weight | role
-----+-----+-----+-----+-----+-----
0        | pghost4  | 1921 | down   | -nan      | standby |
0        | false    | 0    |        |           |
1        | pghost5  | 1921 | up     | -nan      | primary |
```

```
0      | true      | 0
(2 rows)
```

```
pgghost4 down
standby pgghost5 up
primary
```

```
pgpool pgpool
pgpool pgpool
PostgreSQL
```

```
pgghost4 pgpool
pgghost5 pgpool VIP
pgghost4
```

```
pgghost4 pgghost4 pgpool
pgghost5 pgpool

```

```
[root@pgghost4 pgpool]# reboot
```

```
pgghost5 /tmp/failover.log
```

```
Thu Oct 12 11:43:13 CST 2017
falling_node=0
```

```
old_primary=0
new_primary=pghost5
pgdata=/data1/pg10/pg_root
Local promote
```

```
pgghost4pgghost5
/etc/pgpool-II/failover_stream.sh
```

```
pghost5
```

```
[postgres@pghost5 ~]$ pg_controldata | grep cluster
Database cluster state:          in production
```

```
pghost5
```

```
pgpool
```

```
[postgres@pghost5 ~]$ psql -h 192.168.26.72 -p 9999 postgres
postgres=# show pool_nodes;
   node_id | hostname | port | status | lb_weight | role
-----+-----+-----+-----+-----+-----
         0 | pghost4  | 1921 | down   | -nan      |
standby   | 0        | false | 0      |
         1 | pghost5  | 1921 | up     | -nan      |
primary   | 0        | true  | 0      |
(2 rows)
```

```
pghost4down
standbypghost5up
```


primary

pgghost5 /var/log/messages
pgpool

```
...
Oct 12 11:43:15 pgghost5 ntpd[6835]: Listen normally on 9
bond0:1 192.168.26.72 UDP 123
Oct 12 11:43:17 pgghost5 pgpool: watchdog escalation[44067]:
[125-1] 2017-10-12 11:43:17: pid 44067: LOG:  successfully
acquired the delegate IP:"192.168.26.72"
Oct 12 11:43:17 pgghost5 pgpool: watchdog escalation[44067]:
[125-2] 2017-10-12 11:43:17: pid 44067: DETAIL:  'if_up_cmd'
returned with success
Oct 12 11:43:17 pgghost5 pgpool: watchdog[42054]: [145-1]
2017-10-12 11:43:17: pid 42054: LOG:  watchdog escalation
process with pid: 44067 exit with SUCCESS.
...
```

pgpool.conf if_up_cmd
VIP ip addr

```
[postgres@pgghost5 ~]$ ip a
...
12: bond0: <BROADCAST,MULTICAST,MASTER,UP,LOWER_UP> mtu 1500
qdisc noqueue state UP
    link/ether a0:36:9f:9d:c7:5f brd ff:ff:ff:ff:ff:ff
    inet 192.168.26.58/24 brd 192.168.26.255 scope global
bond0
    inet 192.168.26.72/24 scope global secondary bond0:1
...
```

VIP pgghost5

pgpool 12.5

pgpool 12.5

- pgpool pgpool

- PostgreSQL pgpool PostgreSQL

- pgpool PostgreSQL pgpool VIP pgpool

pgpool



pgpool failover_stream.sh Keepalived+

14.1.5 pgpool

pgpool pgpool

pgpool pool_passwd

```
[postgres@pghost4 ~]$ psql -h 192.168.26.72 -p 9999 postgres
pguser
psql: FATAL:  md5 authentication failed
DETAIL:  pool_passwd file does not contain an entry for
"pguser"
```

pgpool enable_pool_hba
pgpool pool_hba.conf
pool_hba.conf md5
pool_passwd 14.1.2
pool_passwd repuser
pool_passwd

repuser:md54f87427f75b5a59ba0abffe11a6f79a8

pgpool

pcp

```
[postgres@pghost5 ~]$ pcp_node_info -h 192.168.26.72 -U
pgpool 0
Password:
```

```
FATAL: authentication failed for user "pgpool"
DETAIL:  username and/or password does not match
```

14.1.3 pgpool 2.2.12 中 pgpool.conf の設定

pgpool.conf の設定

pgpool.conf の設定

```
# USERID:MD5PASSWD
pgpool:ba777e4c2f15c11ea8ac3be7e0440aa0
```

pgpool の VIP の設定

pgpool の VIP の設定

pgpool.conf の設定

if_up_cmd, if_down_cmd の設定

VIP の設定

```
if_up_cmd = 'ip addr add $_IP_$ /24 dev bond0 label bond0:1'
if_down_cmd = 'ip addr del $_IP_$ /24 dev bond0'
```

pgpool の設定

bond0 の設定

VIP の設定

bond0 の設定

pgpool の設定

VIP の設定

pgpool の設定

if_up_cmd, if_down_cmd の設定

ip addr add, ip addr del の設定

./failover_stream.sh
falling_node old_primary
PostgreSQL

PostgreSQL
pcp_attach_node pgpool

pgpool

```
postgres=# show pool_nodes;
 node_id | hostname | port | status | lb_weight | role
| select_cnt | load_balance_node | replication_delay
-----+-----+-----+-----+-----+-----
-+-----+----
    0      | pghost4  | 1921 | up      | -nan      | primary
| 0      | true      | 0
    1      | pghost5  | 1921 | up      | -nan      | standby
| 0      | false     | 0
(2 rows)
```

pghost5

```
[postgres@pghost5 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

pgpool show pool_nodes
pghost5 down pghost5

```
[postgres@pghost5 ~]$ pg_ctl start
```

pgpool

```
postgres=# show pool_nodes;
   node_id | hostname | port | status | lb_weight | role
| select_cnt | load_balance_node | replication_delay
-----+-----+-----+-----+-----+-----
   0      | pghost4  | 1921 | up      | -nan      | primary
| 0         | true      | 0     |
   1      | pghost5  | 1921 | down    | -nan      | standby
| 0         | false     | 0     |
(2 rows)
```

pghost5 down
pcp_attach_node PostgreSQL
pgpool

```
[postgres@pghost4 ~]$ pcp_attach_node -h 192.168.26.72 -U
pgpool 1
Password:
pcp_attach_node -- Command Successful
```

pcp_attach_node

```
pcp_attach_node [options...] [node_id]
```

```
node_id PostgreSQL pgghost4  
0 pgghost5 1 pgghost5  
pgpool
```

```
pgpool show pool_nodes  
pgghost5 up
```


14.2 Keepalived+负载均衡

pgpool PostgreSQL 负载均衡
Keepalived 负载均衡
Keepalived PostgreSQL
Keepalived Keepalived
VIP

PostgreSQL

PostgreSQL10 14-2

14-2 Keepalived+负载均衡

主机名	组件	IP地址	版本	VIP
pghost4	主库	192.168.26.57	PostgreSQL10	192.168.26.72
	Keepalived 主	192.168.26.57	Keepalived-1.3.7	
pghost5	备库	192.168.26.58	PostgreSQL10	
	Keepalived 备	192.168.26.58	Keepalived-1.3.7	

14.2.1 Keepalived+负载均衡

Keepalived+高可用数据库14-2

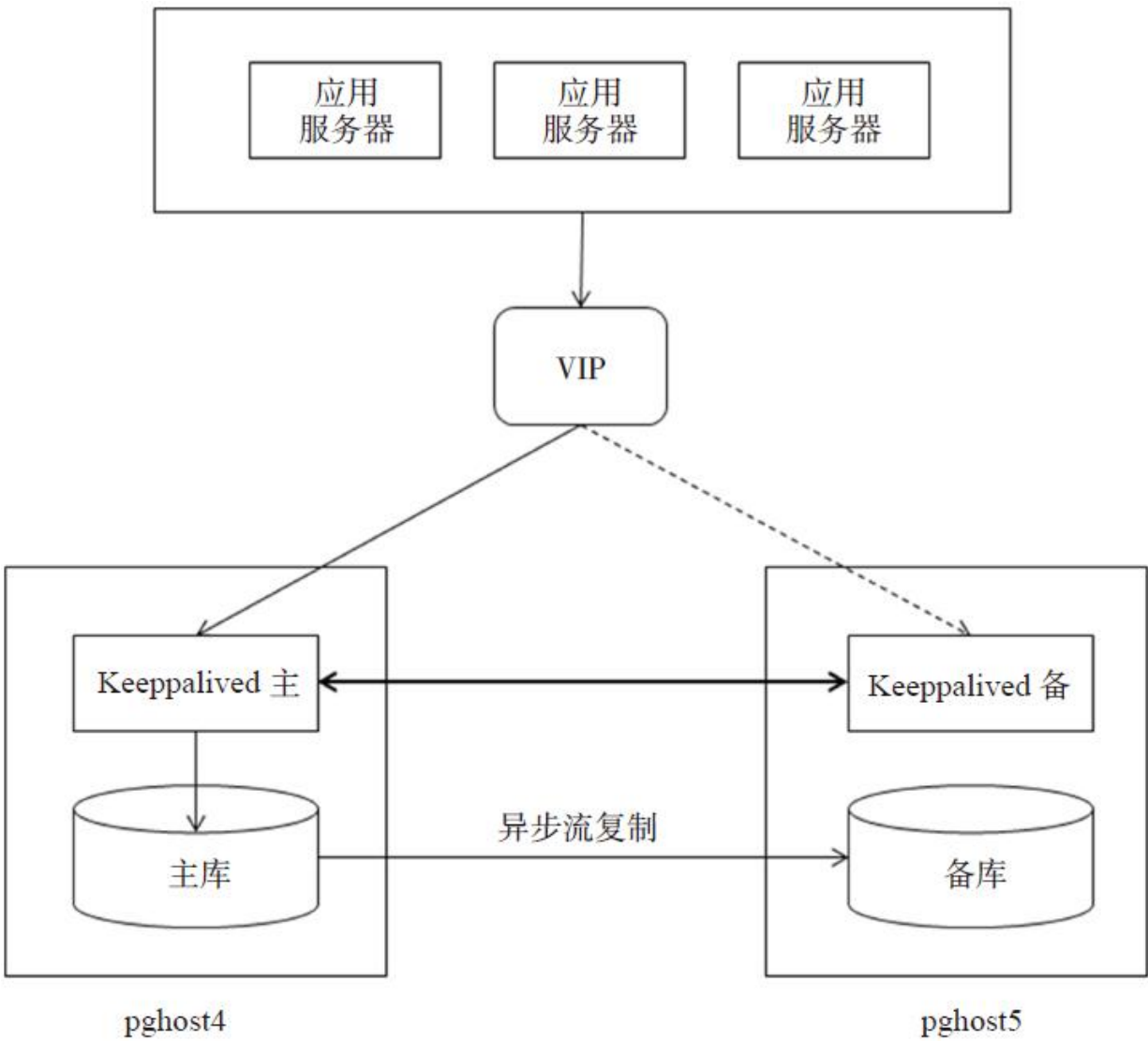


图14-2 Keepalived+高可用数据库

Keepalived配置在pghost4
Keepalived配置在pghost5
数据库配置在pghost4和pghost5

Keepalived
Keepalived
KeepalivedVIPKeepalived

14.2.2 Keepalived+

1
pghost4
pghost5
2KeepalivedKeepalived
sr_delayKeepalived
last_alive

```
postgres=# CREATE ROLE keepalived NOSUPERUSER NOCREATEDB
          login ENCRYPTED PASSWORD 'keepalived';
CREATE ROLE
postgres=# CREATE DATABASE keepalived
          WITH OWNER=keepalived
          TEMPLATE=TEMPLATE0
          ENCODING='UTF8';
CREATE DATABASE

postgres=# \c keepalived keepalived
keepalived=> CREATE TABLE sr_delay(id int4,last_alive
timestamp(0) without time zone);
CREATE TABLE
```

sr_delay

```
CREATE FUNCTION cannt_delete ()
RETURNS trigger
LANGUAGE plpgsql AS $$
BEGIN
RAISE EXCEPTION 'You can not delete!';
END; $$;
```

cannt_delete cannt_truncate

```
keepalived=> CREATE TRIGGER cannt_delete BEFORE DELETE ON
sr_delay
FOR EACH ROW EXECUTE PROCEDURE cannt_delete();
CREATE TRIGGER
keepalived=> CREATE TRIGGER cannt_truncate BEFORE TRUNCATE
ON sr_delay
FOR STATEMENT EXECUTE PROCEDURE cannt_delete();
CREATE TRIGGER
```

sr_delay

```
keepalived=> INSERT INTO sr_delay VALUES(1,now());
INSERT 0 1
```

Keepalived PostgreSQL
Keepalived Keepalived
pg_hba.conf

```
# keepalived
host keepalived keepalived 192.168.26.57/32 md5
host keepalived keepalived 192.168.26.58/32 md5
host keepalived keepalived 192.168.26.72/32 md5
```

pg_ctl reload

3pgghost4pgghostKeepalived

Keepalived

<http://www.keepalived.org/software/keepalived-1.3.7.tar.gz>

```
# yum install openssl openssl-devel popt popt-devel
```

pgghost4pgghostKeepalived

```
# tar xvf keepalived-1.3.7.tar.gz
# ./configure --prefix=/usr/local/keepalived
# make
# make install
```

Keepalived

```
# ln -s /usr/local/keepalived/sbin/keepalived /usr/sbin/  
# cp /opt/soft_bak/keepalived-  
1.3.7/keepalived/etc/init.d/keepalived /etc/init.d/  
# cp /usr/local/keepalived/etc/sysconfig/keepalived  
/etc/sysconfig
```

Keepalived
/etc/init.d

```
# service keepalived status  
keepalived is stopped
```

Keepalived Keepalived

```
# chkconfig keepalived on
```

14.2.3 Keepalived

Keepalived

Keepalived

```
# mkdir -p /etc/keepalived
```

##/etc/keepalived/keepalived.conf##

```
! Configuration File for keepalived

global_defs {
    notification_email {
        francs3@163.com
    }
    smtp_server 127.0.0.1
    smtp_connect_timeout 30
    router_id DB1_PG_HA
}

vrrp_script check_pg_alived {
    script "/usr/local/bin/pg_monitor.sh"
    interval 10
    fall 3      # require 3 failures for K0
}

vrrp_instance VI_1 {
    state BACKUP
    nopreempt
    interface bond0
    virtual_router_id 10
    priority 100
    advert_int 1
    authentication {
        auth_type PASS
        auth_pass t9rveMP0Z9S1
    }
    track_script {
        check_pg_alived
    }
    virtual_ipaddress {
        192.168.26.72
    }
    smtp_alert
    notify_master /usr/local/bin/active_standby.sh
}
```

Keepalived Keepalived
priority 90

·global_defs
Keepalived

·vrrp_script 10
/usr/local/bin/pg_monitor.sh
fall 3

·vrrp_instance vrrp
ID backup
nopreempt VIP priority
Keepalived priority
Keepalived Keepalived VIP
bond0

·smtp_alert notify_master
Keepalived
/usr/local/bin/active_standby.sh

/usr/local/bin/pg_monitor.sh

```
#!/bin/bash
#
export PGPORT=1921
export PGUSER=keepalived
```



```

export PGDBNAME=keepalived
export PGDATA=/data1/pg10/pg_root
export LANG=en_US.utf8
export PGHOME=/opt/pgsql
export
LD_LIBRARY_PATH=$PGHOME/lib:/lib64:/usr/lib64:/usr/local/lib
64:/lib:/usr/lib:/usr/local/lib
export PATH=$PGHOME/bin:$PGPOOL_HOME/bin:$PATH:.

```

```

MONITOR_LOG="/tmp/pg_monitor.log"
SQL1="UPDATE sr_delay SET last_alive = now();"
SQL2='SELECT 1;'
# 检查是否处于恢复状态
standby_flg=`psql -p $PGPORT -U postgres -At -c "SELECT
pg_is_in_recovery();"
if [ ${standby_flg} == 't' ]; then
    echo -e "`date +%F\ %T`: This is a standby database,
exit!\n" >> $MONITOR_LOG
    exit 0
fi

```

```

# 检查sr_delay表
echo $SQL1 | psql -At -p $PGPORT -U $PGUSER -d $PGDBNAME >>
$MONITOR_LOG

```

```

# 检查主数据库健康状态
echo $SQL2 | psql -At -h -p $PGPORT -U $PGUSER -d $PGDBNAME
if [ $? -eq 0 ]; then
    echo -e "`date +%F\ %T`: Primary db is health." >>
$MONITOR_LOG
    exit 0
else
    echo -e "`date +%F\ %T`: Attention: Primary db is not
health!" >> $MONITOR_LOG
    exit 1
fi

```

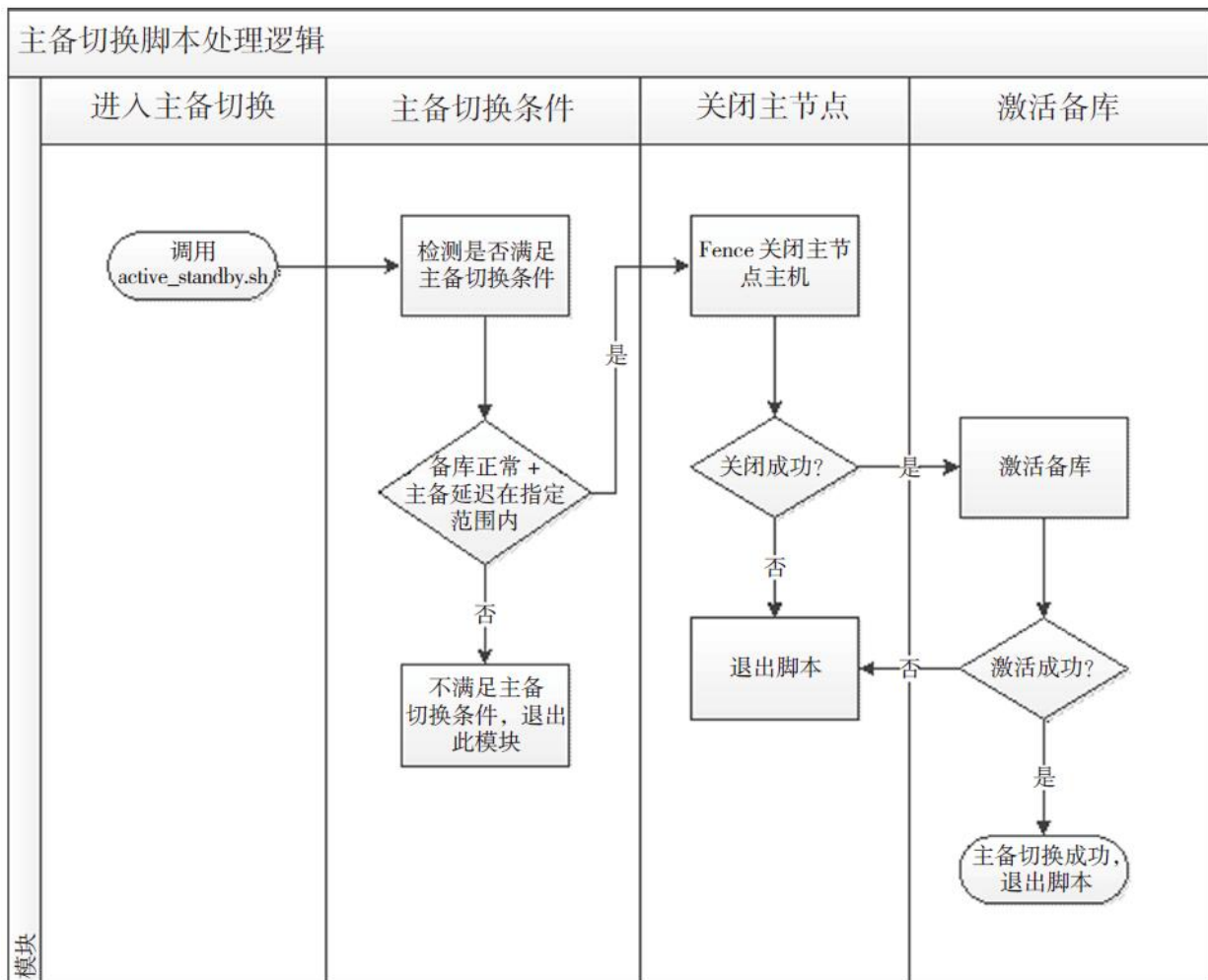
10keepalived.conf
 interval

.

· sr_delay last_alive

Keepalived Keepalived
 Keepalived
 notify_master
 /usr/local/bin/active_standby.sh

/usr/local/bin/active_standby.sh
 14-3



14-3 active_stadnby.sh

active_standby.sh

1

2

1

2

/usr/local/bin/active_standby.sh

```
#!/bin/bash
#
export PGPOR=1921
export PGUSER=keepalived
export PG_OS_USER=postgres
export PGDBNAME=keepalived
export PGDATA=/data1/pg10/pg_root
export LANG=en_US.utf8
export PGHOME=/opt/pgsql
export
LD_LIBRARY_PATH=$PGHOME/lib:/lib64:/usr/lib64:/usr/local/lib
64:/lib:/usr/lib:/usr/local/lib
```

```
export
PATH=/opt/pgbouncer/bin:$PGHOME/bin:$PGPOOL_HOME/bin:$PATH:.
```

```
# 配置, LAG_MINUTES 配置
LAG_MINUTES=60
HOST_IP=`hostname -i`
NOTICE_EMAIL="francs3@163.com"
FAILOVE_LOG='/tmp/pg_failover.log'
```

```
SQL1="SELECT 'this_is_standby' AS cluster_role FROM ( SELECT
pg_is_in_recovery() AS std ) t WHERE t.std is true;"
SQL2="SELECT 'standby_in_allowed_lag' AS cluster_lag FROM
sr_delay WHERE now()-last_alive < interval '$LAG_MINUTES
SECONDS';"
```

```
# 配置IP
FENCE_IP=50.1.225.101
FENCE_USER=root
FENCE_PWD=xxxx
```

```
# VIP 配置
echo -e "`date +%F\ %T`: keepalived VIP switchover!" >>
$FAILOVE_LOG
```

```
# VIP 配置
#echo -e "`date +%F\ %T`: ${HOST_IP}/${PGPORT} VIP 配置
\n\nAuthor: francs(DBA)" | mutt -s "Error: VIP 配置 "
${NOTICE_EMAIL}
```

```
# pg_failover 配置
pg_failover()
{
# FENCE_STATUS 配置1 配置0 配置
# PROMOTE_STATUS 配置1 配置0 配置
FENCE_STATUS=1
PROMOTE_STATUS=1
```

```
# 配置
for ((k=0;k<10;k++))
do
# ipmitool配置X86配置
ipmitool -I lanplus -L OPERATOR -H $FENCE_IP -U
$FENCE_USER -P $FENCE_PWD power reset
if [ $? -eq 0 ]; then
echo -e "`date +%F\ %T`: fence primary db host
```

```

success."
        FENCE_STATUS=0
        break
    fi
sleep 1
done

if [ $FENCE_STATUS -ne 0 ]; then
    echo -e "`date +%F\ %T`: fence failed. Standby will not
promote, please fix it manually."
    return $FENCE_STATUS
fi

# 检查
su - $PG_OS_USER -c "pg_ctl promote"
if [ $? -eq 0 ]; then
    echo -e "`date +%F\ %T`: `hostname` promote standby
success. "
    PROMOTE_STATUS=0
fi

if [ $PROMOTE_STATUS -ne 0 ]; then
    echo -e "`date +%F\ %T`: promote standby failed."
    return $PROMOTE_STATUS
fi

    echo -e "`date +%F\ %T`: pg_failover() function call
success."
    return 0
}

# 检查 standby 状态
# 检查 standby 是否处于 standby 状态
STANDBY_CNT=`echo $SQL1 | psql -At -p $PGPORT -U $PGUSER -d
$PGDBNAME -f - | grep -c this_is_standby`
echo -e "STANDBY_CNT: $STANDBY_CNT" >> $FAILOVER_LOG

if [ $STANDBY_CNT -ne 1 ]; then
    echo -e "`date +%F\ %T`: `hostname` is not standby
database failover not allowed! " >> $FAILOVER_LOG
    exit 1
fi

# 检查 standby 是否处于 standby 状态, LAG=1 检查 standby 是否处于 standby
LAG=`echo $SQL2 | psql -At -p $PGPORT -U $PGUSER -d

```

```

$PGDBNAME | grep -c standby_in_allowed_lag`
echo -e "LAG: $LAG" >> $FAILOVE_LOG

if [ $LAG -ne 1 ]; then
    echo -e "`date +%F\ %T`: `hostname` is laged far
$LAG MINUTES SECONDS from primary , failover not allowed! "
>> $FAILOVE_LOG
    exit 1
fi

# 1. standby_in_allowed_lag 2. standby_in_allowed_lag
if [ $STANDBY_CNT -eq 1 ] && [ $LAG -eq 1 ]; then
    pg_failover >> $FAILOVE_LOG
    if [ $? -ne 0 ]; then
        echo -e "`date +%F\ %T`: pg_failover failed." >>
$FAILOVE_LOG
        exit 1
    fi
fi

# 1. pg_failover
# 2. pg_failover

# 1. pg_failover
# 2. pg_failover

```

```

FENCE_IP FENCE_USER
FENCE_PWD

```

```


```

```

# chmod 700 /usr/local/bin/pg_monitor.sh
# chmod 700 /usr/local/bin/active_standby.sh

```



1. 在 /etc/crontab 文件中添加如下内容
 0 * * * * root pg_monitor.sh
 0 * * * * root active_standby.sh

14.2.4 Keepalived 配置

Keepalived 是一个高可用解决方案，它结合了 VRRP 和 LVS。

· Keepalived 配置文件 keepalived.conf

· Keepalived 配置文件 keepalived.conf

· Keepalived 配置文件 keepalived.conf

Keepalived 配置文件 keepalived.conf

pghost4 Keepalived
 pghost5 Keepalived VIP
 pghost4

pghost4 Keepalived
Keepalived

```
[root@pghost4 ~]# ps -ef | grep keepalived | grep -v grep
root      7527      1  0 10:31 ?        00:00:00 keepalived -
D
root      7528    7527  0 10:31 ?        00:00:00 keepalived -
D
root      7529    7527  0 10:31 ?        00:00:00 keepalived -
D
[root@pghost4 ~]# kill 7527
```

pghost5 /var/log/messages
[]

```
Oct 16 11:00:55 pghost5 Keepalived_vrrp[6561]:
VRRP_Instance(VI_1) Transition to MASTER STATE
Oct 16 11:00:56 pghost5 Keepalived_vrrp[6561]:
VRRP_Instance(VI_1) Entering MASTER STATE
Oct 16 11:00:56 pghost5 Keepalived_vrrp[6561]:
VRRP_Instance(VI_1) setting protocol VIPs.
...
Oct 16 11:00:57 pghost5 ntpd[6917]: Listen normally on 6
bond0 192.168.26.72 UDP 123
Oct 16 11:01:01 pghost5 Keepalived_vrrp[6561]: Sending
gratuitous ARP on bond0 for 192.168.26.72
...
```

pghost5 Keepalived
VIP 192.168.26.72

pghost5 ip addr

```
[root@pghost5 ~]# ip addr
...
12: bond0: <BROADCAST,MULTICAST,MASTER,UP,LOWER_UP> mtu 1500
qdisc noqueue state UP
    link/ether a0:36:9f:9d:c7:5f brd ff:ff:ff:ff:ff:ff
    inet 192.168.26.58/24 brd 192.168.26.255 scope global
bond0
    inet 192.168.26.72/32 scope global bond0
...
```

[[VIP]]pghost5[[

pghost5[[/tmp/pg_failover.log[[
[[active_standby.sh[[

```
2017-10-16 11:00:56: keepalived VIP switchover!
STANDBY_CNT: 1
LAG: 1
Chassis Power Control: Reset
2017-10-16 11:00:56: fence primary db host success.
waiting for server to promote.... done
server promoted
2017-10-16 11:00:56: pghost5 promote standby success.
2017-10-16 11:00:56: pg_failover() function call success.
```

[[Keepalived[[VIP[[
[[pghost5[[pghost4[[
[[

[[pghost5[[

```
[postgres@pghost5 pg_root]$ pg_controldata | grep cluster
Database cluster state:                in production
```

pgghost5

pgghost4 Keepalived
pgghost5 Keepalived
VIP pgghost4 pgghost4
\$PGDATA recovery.done
recovery.conf

12.5

pgghost4

```
[postgres@pgghost4 ~]$ pg_ctl stop -m fast
waiting for server to shut down.... done
server stopped
```

pgghost4

```
...
Oct 16 11:20:35 pgghost4 Keepalived_vrrp[7436]:
/usr/local/bin/pg_monitor.sh exited with status 1
Oct 16 11:20:35 pgghost4 Keepalived_vrrp[7436]:
VRRP_Script(check_pg_alived) failed
```

```
Oct 16 11:20:35 pghost4 Keepalived_vrrp[7436]:  
VRRP_Instance(VI_1) Entering FAULT STATE  
Oct 16 11:20:35 pghost4 Keepalived_vrrp[7436]:  
VRRP_Instance(VI_1) removing protocol VIPs.  
Oct 16 11:20:35 pghost4 Keepalived_vrrp[7436]:  
VRRP_Instance(VI_1) Now in FAULT state  
Oct 16 11:20:36 pghost4 ntpd[6888]: Deleting interface #6  
bond0, 192.168.26.72#123, interface stats: received=0,  
sent=0, dropped=0, active_time=442 secs  
...
```

```
#####check_pg_alived#####  
Keepalived#####VIP#####pghost4#####bond0#####  
#####pghost4#####pghost5#####  
Keepalived#####active_standby.sh  
#####
```

```
#####pghost5#####
```

```
...  
Oct 16 11:20:38 pghost5 ntpd[6917]: Listen normally on 6  
bond0 192.168.26.72 UDP 123  
Oct 16 11:20:41 pghost5 Keepalived_vrrp[6561]: Sending  
gratuitous ARP on bond0 for 192.168.26.72  
Oct 16 11:20:41 pghost5 Keepalived_vrrp[6561]:  
VRRP_Instance(VI_1) Sending/queueing gratuitous ARPs on  
bond0 for 192.168.26.72  
...
```

```
#####pghost5#####Keepalived#####  
#####VIP 192.168.26.72#####
```

pgghost5/tmp/pg_failover.log

```
2017-10-16 11:20:36: keepalived VIP switchover!
STANDBY_CNT: 1
LAG: 1
Chassis Power Control: Reset
2017-10-16 11:20:37: fence primary db host success.
waiting for server to promote.... done
server promoted
2017-10-16 11:20:37: pgghost5 promote standby success.
2017-10-16 11:20:37: pg_failover() function call success.
```

Keepalived VIP

pgghost5

```
[postgres@pgghost5 pg_root]$ pg_controldata | grep cluster
Database cluster state:          in production
```

pgghost5

pgghost5
ipmitool power reset pgghost5

pgghost4 Keepalived
pgghost5

· Keepalived Keepalived
PostgreSQL

· Keepalived
PostgreSQL

· Keepalived
PostgreSQL

Keepalived
pgpool+
pgpool

14.3 ☐☐☐☐

```
pgpool
Keepalived
pgpool
pgpool
pgpool
Keepalived
```

PostgreSQL

15

PostgreSQL
BUG
PostgreSQL
PostgreSQL
PostgreSQL
PostgreSQL

15.1 □□□□

PostgreSQL 10
PostgreSQL
9.6.14
9.6.14
2017
10
PostgreSQL 10
PostgreSQL
10.0
PostgreSQL
2
5
8
11
BUG

PostgreSQL Long-term support LTS 5

PostgreSQL

版 本	发 行 时 间	特 性
8.4	2009-07-01	窗口函数，列级权限，并行数据库恢复，公用表表达式和递归查询
9.0	2010-09-20	内置流复制、热备、支持 64 位 Windows 操作系统
9.1	2011-09-12	同步流复制，UNLOGGED 表，可串行化快照隔离，可写公用表表达式，SELinux 集成，外部扩展，外部表
9.2	2012-09-10	级联流复制，只用索引的扫描，原生 json 支持，范围类型，pg_receivexlog 工具，Space-Partitioned GiST 索引
9.3	2013-09-09	自定义后台工作进程，数据校验，专用 JSON 运算符，LATERAL JOIN，更快的 pg_dump, pg_isready 服务器监控工具，触发器，视图，可写外部表，物化视图
9.4	2014-12-18	JSONB 数据类型，用于更改配置值的 ALTER SYSTEM 语句，后台工作进程动态注册 / 启动 / 停止，逻辑解析 API，Linux 大页支持，pg_prewarm 缓存预热
9.5	2016-01-07	UPSERT, 行级安全, 数据抽样, BRIN 索引
9.6	2016-09-29	并行查询，PostgreSQL 外部表增加排序合并操作下推、多个同步复制从库、vacuum 大表速度增加
10	2017-10-05	逻辑复制，原生分区表、并行查询增强

15.2 □□□□□

PostgreSQL

9.6.4 9.6.5

□ □ □ □ □ □ □ □ □ □ □

```
psql
SELECT version
```

```
[postgres@pghost1 ~]$ /usr/pgsql-9.6/bin/psql -p 1921 -U
pguser mydb
psql (9.6.4)
Type "help" for help.
mydb=> SELECT version();
        version
```

```
PostgreSQL 9.6.4 on x86_64-pc-linux-gnu, compiled by gcc
(GCC) 4.4.7 20120313 (Red Hat 4.4.7-18), 64-bit
(1 row)
```

--	--	--	--	--	--	--	--

psql (9.6.5, server 9.6.4)
Type "help" for help.

mydb=> SELECT version();

version

```
[postgres@pghost1 ~]$ /usr/pgsql-9.6/bin/psql -U pguser -p
1921 mydb
psql (9.6.5, server 9.6.4)
Type "help" for help.
mydb=> SELECT version();
      version
```

```
-----
PostgreSQL 9.6.5 on x86_64-pc-linux-gnu, compiled by gcc
(GCC) 4.4.7 20120313 (Red Hat 4.4.7-18), 64-bit
(1 row)
```

psql (9.6.5, server 9.6.4)
Type "help" for help.

15.3 备份与恢复

PostgreSQL 数据库的备份与恢复是数据库管理员的重要任务之一。本章将介绍 PostgreSQL 10 的备份与恢复方法，包括物理备份、逻辑备份以及 pg_upgrade 和 pg_logical 等工具的使用。

本章将介绍 PostgreSQL 10 的备份与恢复方法，包括物理备份、逻辑备份以及 pg_upgrade 和 pg_logical 等工具的使用。本章将介绍 PostgreSQL 10 的备份与恢复方法，包括物理备份、逻辑备份以及 pg_upgrade 和 pg_logical 等工具的使用。

15.3.1 使用 pg_dumpall 进行逻辑备份

pg_dumpall 是一个用于导出整个 PostgreSQL 数据库的逻辑备份工具。它可以将数据库中的所有数据、表结构、索引、视图、序列、触发器、存储过程和函数等导出到一个文本文件中。该工具支持多种输出格式，包括文本、SQL 和 XML 等。

使用 pg_dumpall 进行逻辑备份的语法如下：

PostgreSQL 9.3 及以上版本支持 pg_dumpall 工具。以下是一个使用 pg_dumpall 进行逻辑备份的示例：

1. 备份数据库

rpm --prefix /usr/pgsql-10 --install postgresql10

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/initdb -D /pgdata/10/data/
```

2. 复制配置文件

使用 `pg_dumpall` 命令备份数据库，并复制配置文件 `pg_hba.conf`、`pg_ident.conf`、`postgresql.conf` 到新的数据目录。

复制配置文件 `pg_hba.conf`、`pg_ident.conf` 和 `postgresql.conf`。

```
[postgres@pghost1 ~]$ cp /pgdata/9.3/data/pg_hba.conf /pgdata/10/data/
[postgres@pghost1 ~]$ cp /pgdata/9.3/data/pg_ident.conf /pgdata/10/data
```

复制 `postgresql.conf` 文件，并添加 `include 'pg93.conf'` 语句，以包含旧版本的配置文件。

```
[postgres@pghost1 ~]$ cp /pgdata/9.3/data/postgresql.conf /pgdata/10/data/pg93.conf
```

3. 启动 PostgreSQL

在启动 PostgreSQL 之前，需要先检查配置文件的正确性。如果配置文件存在错误，PostgreSQL 将无法启动。

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_ctl -D /pgdata/10/data/
start
waiting for server to start....
2017-10-06 06:51:12.950 GMT [652] LOG:  unrecognized configuration
parameter "checkpoint_segments" in file
"/pgdata/10/data/postgresql.auto.conf" line 28
2017-10-06 06:51:12.950 GMT [652] FATAL:  configuration file
"/pgdata/10/data/postgresql.auto.conf" contains errors
        stopped waiting
pg_ctl: could not start server
```

在 PostgreSQL 9.3 到 10 的发布说明中，提到了关于配置文件的更改。请仔细阅读这些说明，以确保配置文件的正确性。

4. 配置 PostgreSQL

在配置 PostgreSQL 时，需要修改配置文件。例如，在 `pg_hba.conf` 文件中，可以配置数据库的访问权限。在 `pg_dumpall` 命令中，可以指定使用 `iptables` 防火墙规则。

在 `postgresql.conf` 文件中，可以配置数据库的默认事务只读模式。例如，可以将 `default_transaction_read_only` 设置为 `on`。

在 PostgreSQL 9.4 中，可以使用 `ALTER DATABASE` 命令来修改数据库的属性。

```
mydb=# ALTER DATABASE mydb SET default_transaction_read_only = TRUE;
```

```

pgbouncer
pgbouncer

```

```
[postgres@pghost1 ~]$ /usr/pgsql-9.3/bin/psql -p 5433 -U pguser mydb
psql (9.3.15, server 9.3.15)
Type "help" for help.
mydb=> SELECT state,COUNT(*) FROM pg_stat_activity WHERE pid <>
pg_backend_pid() GROUP BY state;
   state | count
-----+-----
(0 rows)
```

5.

```

pg_dumpall
pg_dumpall

```

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_dumpall -p 5432 >
backup.sql
```

```


```

```
/usr/pgsql-10/bin/psql -p 1921 -f backup.sql
```

```

Linux

```

```

[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_dumpall -p 5432 |
/usr/pgsql-10/bin/psql -p 1921
SET
...
CREATE ROLE
ALTER ROLE
ERROR:  role "postgres" already exists
ALTER ROLE
CREATE DATABASE
REVOKE
GRANT
You are now connected to database "mydb" as user "postgres".
...
SET
CREATE EXTENSION
COMMENT
...
ALTER TABLE
COPY 3
  setval
-----
  3
(1 row)
...

```

```

                                □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□“ERROR□role"postgres"already exists”□□
postgres□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```

15.3.2 □□pg_upgrade□□□□□□□□

```

                                □□□□□□□□□□pg_upgrade□□□□pg_upgrade□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```


PostgreSQL 9.3から10にアップグレードする
pg_upgradeの使い方

1.pg_upgrade

pg_upgradeは、PostgreSQL 9.3から10にアップグレードするためのツールです。PostgreSQL 9.3から10にアップグレードする際には、pg_upgradeを使用する必要があります。

pg_upgradeの使い方

-b, --old-bindir=BINDIR	旧PostgreSQLのインストールディレクトリ
-B, --new-bindir=BINDIR	新PostgreSQLのインストールディレクトリ
-c, --check	アップグレード前にチェックを行う
-d, --old-datadir=DATADIR	旧PostgreSQLのデータディレクトリ
-D, --new-datadir=DATADIR	新PostgreSQLのデータディレクトリ
-j, --jobs	並行して実行するジョブの数。CPUの個数に設定すると最適です。
-k, --link	シンリンクを使用する
-o, --old-options=OPTIONS	旧PostgreSQLの起動オプション
-O, --new-options=OPTIONS	新PostgreSQLの起動オプション
-p, --old-port=PORT	旧PostgreSQLのポート番号。デフォルトは5432。
-P, --new-port=PORT	新PostgreSQLのポート番号。デフォルトは5432。
-r, --retain	旧PostgreSQLのデータディレクトリを保持する

pg_upgrade -c -b /usr/pgsql-9.3/bin -B /usr/pgsql-10/bin -d /usr/pgsql-9.3/data -D /usr/pgsql-10/data

pg_upgrade -k Linkは、シンリンクを使用するオプションです。シンリンクを使用すると、旧PostgreSQLのデータディレクトリを保持することができます。

pg_upgrade link 是硬链接，所以升级的时候，pg_upgrade 会先检查旧版本的数据目录和二进制目录是否一致，如果一致，则不会创建新的数据目录和二进制目录，而是直接链接到旧版本的数据目录和二进制目录。

2. pg_upgrade

1. 检查 PostgreSQL 版本是否一致

2. 检查数据目录和二进制目录是否一致

3. 检查数据目录和二进制目录是否一致

pg_upgrade 命令支持 --check 选项，用于检查旧版本的数据目录和二进制目录是否一致。如果检查通过，则不会创建新的数据目录和二进制目录，而是直接链接到旧版本的数据目录和二进制目录。

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_upgrade -b /usr/pgsql-9.3/bin -B /usr/pgsql-10/bin -d /pgdata/9.3/data/ -D /pgdata/10/data/ -k -c
Performing Consistency Checks
-----
Checking cluster versions                                ok
...
Checking for prepared transactions                       ok
*Clusters are compatible*
```

如果检查通过，则输出 "Clusters are compatible"。如果检查失败，则输出 "Failure: exiting"。

Your installation references loadable libraries that are missing from the new installation. You can add these libraries to the new installation, or remove the functions using them from the old installation. A list of problem libraries is in the file:
loadable_libraries.txt

loadable_libraries.txt

```
[postgres@pghost1 ~]$ cat loadable_libraries.txt
could not load library "$libdir/postgis-2.3": ERROR:  could not
access file "$libdir/postgis-2.3": No such file or directory
could not load library "$libdir/rtpostgis-2.3": ERROR:  could not
access file "$libdir/rtpostgis-2.3": No such file or directory
```

warning

```
Checking for hash indexes warning
Your installation contains hash indexes.  These indexes have
different
internal formats between your old and new clusters, so they must be
reindexed with the REINDEX command.  After upgrading, you will be
given
REINDEX instructions.
```

PostgreSQL 10 hash

```
mydb=# \d+ tbl
```

```
...
```

```
Indexes:
```

```
    "tbl_hash_idx" hash (column_name) INVALID
```

warning hash

```
mydb=# CREATE INDEX CONCURRENTLY ON th USING HASH(column_name);
CREATE INDEX
```

```
mydb=# DROP INDEX tbl_hash_idx;  
DROP INDEX
```

4 pg_upgrade

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_upgrade -b /usr/pgsql-  
9.3/bin -B /usr/pgsql-10/bin -d /pgdata/9.3/data/ -D /pgdata/10/data/  
Performing Consistency Checks
```

```
-----  
Checking cluster versions                                ok  
Checking database user is the install user              ok  
Checking database connection settings                    ok  
Checking for prepared transactions                       ok  
Checking for reg* data types in user tables              ok  
Checking for contrib/iso with bigint-passing mismatch  ok  
Checking for invalid "unknown" user columns             ok  
Checking for roles starting with "pg_"                  ok  
Checking for incompatible "line" data type              ok  
Creating dump of global objects                          ok  
Creating dump of database schemas
```

```
-----  
Checking for presence of required libraries              ok  
Checking database user is the install user              ok  
Checking for prepared transactions                       ok  
If pg_upgrade fails after this point, you must re-initdb the  
new cluster before continuing.  
Performing Upgrade
```

```
-----  
Analyzing all rows in the new cluster                    ok  
Freezing all rows in the new cluster                     ok  
Deleting files from new pg_xact                          ok  
Copying old pg_clog to new server                        ok  
Setting next transaction ID and epoch for new cluster    ok  
Deleting files from new pg_multixact/offsets             ok  
Copying old pg_multixact/offsets to new server           ok  
Deleting files from new pg_multixact/members             ok  
Copying old pg_multixact/members to new server           ok  
Setting next multixact ID and offset for new cluster     ok  
Resetting WAL archives                                   ok  
Setting frozenxid and minmxid counters in new cluster    ok  
Restoring global objects in the new cluster              ok  
Restoring database schemas in the new cluster
```

```
-----  
Copying user relation files                              ok  
-----  
Setting next OID for new cluster                         ok  
Sync data directory to disk                             ok  
Creating script to analyze new cluster                   ok
```

```
Creating script to delete old cluster          ok
Checking for hash indexes                     ok
Upgrade Complete
-----
Optimizer statistics are not transferred by pg_upgrade so,
once you start the new server, consider running:
    ./analyze_new_cluster.sh
Running this script will delete the old cluster's data files:
    ./delete_old_cluster.sh
```

pg_upgrade
“Upgrade Complete”

5pg_upgradelink

link
Extension

```
mydb=# \db
          List of tablespaces
   Name   | Owner   | Location
-----+-----+-----
 pg_default | postgres |
 pg_global  | postgres |
 pgtablespace | postgres | /pgdata/pgtablespace
(3 rows)
mydb=# \dx
                                List of installed extensions
   Name   | Version | Schema  | Description
-----+-----+-----+-----
 pg_stat_statements | 1.1     | public  | track execution
 statistics of all SQL statements executed
 plpgsql   | 1.0     | pg_catalog | PL/pgSQL procedural
 language
(2 rows)
```

pg_upgrade

```

[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_upgrade -b /usr/pgsql-
9.3/bin -B /usr/pgsql-10/bin -d /pgdata/9.3/data/ -D /pgdata/10/data/
-k
Performing Consistency Checks
-----
Checking cluster versions                                ok
...
Checking for prepared transactions                      ok
If pg_upgrade fails after this point, you must re-initdb the
new cluster before continuing.
Performing Upgrade
-----
Analyzing all rows in the new cluster                   ok
Freezing all rows in the new cluster                   ok
Deleting files from new pg_xact                        ok
Copying old pg_clog to new server                      ok
Setting next transaction ID and epoch for new cluster  ok
Deleting files from new pg_multixact/offsets           ok
Copying old pg_multixact/offsets to new server         ok
Deleting files from new pg_multixact/members          ok
Copying old pg_multixact/members to new server        ok
Setting next multixact ID and offset for new cluster   ok
Resetting WAL archives                                ok
Setting frozenxid and minmxid counters in new cluster  ok
Restoring global objects in the new cluster            ok
Restoring database schemas in the new cluster          ok
                                                         ok
Adding ".old" suffix to old global/pg_control          ok
If you want to start the old cluster, you will need to remove
the ".old" suffix from /pgdata/9.3/data/global/pg_control.old.
Because "link" mode was used, the old cluster cannot be safely
started once the new cluster has been started.
Linking user relation files                             ok
                                                         ok
Setting next OID for new cluster                       ok
Sync data directory to disk                           ok
Creating script to analyze new cluster                 ok
Creating script to delete old cluster                  ok
Checking for hash indexes                             ok
Upgrade Complete
-----
Optimizer statistics are not transferred by pg_upgrade so,
once you start the new server, consider running:
    ./analyze_new_cluster.sh
Running this script will delete the old cluster's data files:
    ./delete_old_cluster.sh

```

```
pg_upgrade pg_upgrade
pg_control pg_control.old
pg_control.old pg_control

```

6

```
pg_upgrade

```

```
pg_upgrade
analyze_new_cluster.sh
vacuumdb--analyze-in-stages VACUUM
VACUUM

```

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/vacuumdb -a --analyze-in-stages -h pghost1 -p 1921
```

7

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_ctl -D /pgdata/10/data/start
```

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/psql -p 1921 mydb
psql (10.0)
Type "help" for help.
mydb=# SELECT version();
        version
```

```
-----
PostgreSQL 10.0 on x86_64-pc-linux-gnu, compiled by gcc (GCC) 4.4.7
20120313 (Red Hat 4.4.7-18), 64-bit
(1 row)
```

8. 检查配置

检查配置是否正确，确保所有配置项都已正确设置。

3. 配置pg_upgrade

配置pg_upgrade，确保所有配置项都已正确设置。

1. 检查配置

PostgreSQL 配置项检查，确保所有配置项都已正确设置。

```
postgres=# SELECT application_name,client_addr,sync_state FROM
pg_stat_replication;
 application_name | client_addr | sync_state
-----+-----+-----
 walreceiver      | 10.191.136.3 | async
(1 row)
```

检查sync_state是否为async，如果是sync，则需要等待。

配置postgresql.conf，设置"synchronous_standby_names"为所有 standby 节点。

执行 reload 命令，使配置生效。

確認データベースシステム識別子の一致

“Database system identifier”を確認

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_controldata
/pgdata/10/data_master | grep "Database system identifier"
Database system identifier:          6504563477800205659
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pg_controldata
/pgdata/10/data_slave | grep "Database system identifier"
Database system identifier:          6504563477800205659
```

pg_upgradeを実行

pg_upgradeを実行

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_upgrade -b /usr/pgsql-
9.3/bin -B /usr/pgsql-10/bin -d /pgdata/9.3/data_master/ -D
/pgdata/10/data_master/ -k
```

バックアップ

4バックアップ

バックアップ

```
[postgres@pghost2 ~]$ cp /pgdata/9.3/data_slave/*.conf
/pgdata/9.3/backup/
```

pg_upgradeを実行

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pg_upgrade -b /usr/pgsql-
9.3/bin -B /usr/pgsql-10/bin -d /pgdata/9.3/data_slave/ -D
/pgdata/10/data_slave/ -k
```

数据库系统标识符“Database system identifier”

```
2017-11-28 18:42:24.692 CST,,,82112,,5a44ca90.140c0,1,,2017-11-28
18:42:24 CST,,0,FATAL,XX000,"database system identifier differs
between the primary and standby","The primary's identifier is
6504542580012252832,
the standby's identifier is
6504544106935856722.",,,,,,,"WalReceiverMain, walreceiver.c:347",,""
```

5

Master

```
[postgres@pghost1 ~]$ /usr/pgsql-10/bin/pg_ctl -D
/pgdata/10/data_master start
```

6

recovery.conf

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pg_ctl -D
/pgdata/10/data_slave start
```

7

```
mydb=# select version();
version
```

```
-----
PostgreSQL 10.1 on x86_64-pc-linux-gnu, compiled by gcc (GCC)
4.4.7 20120313 (Red Hat 4.4.7-18), 64-bit
(1 row)
```

```


```

```
mydb=# SELECT client_addr, sync_state FROM pg_stat_replication;
 client_addr | sync_state
-----+-----
 pghost2    | sync
(1 row)
```

15.3.3 pglogical

1.pglogical

2ndQuadrant pglogical Extension
 PostgreSQL
 pglogical Node PostgreSQL
 Providers Subscribers Node Replication
 Set

2.pglogical

```


```

```
pglogical
```

·PostgreSQL 9.4

- PRIMARY KEY REPLICATION IDENTITY

- UNLOGGED TEMPORARY

- Database

- DDL DDL

pglogical.replicate_ddl_command

pglogical
pglogical

3. pglogical

PostgreSQL 9.4 PostgreSQL 10
pglogical

1

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/initdb -k -D /pgdata/10/data_master
```

2

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pg_ctl -D /pgdata/10/data_master start
```

3 schema

```
[postgres@pghost2 ~]$ /usr/pgsql-10/bin/pg_dump -Fp -v -c -s -h  
pghost1 -p 1921 -d mydb |  
/usr/pgsql-10/bin/psql -h pghost2 -p 1922 -d mydb
```

4 pglogical

pglogical repo RPM

```
[root@pghost1 ~]# yum install  
http://packages.2ndquadrant.com/pglogical/yum-repo-rpms/pglogical-  
rhel-1.0-3.noarch.rpm
```

pglogical

```
# PostgreSQL 9.4  
[root@pghost1 ~]# yum install postgresql94-pglogical.x86_64  
# PostgreSQL 10  
[root@pghost2 ~]# yum install postgresql10-pglogical.x86_64
```

5 pglogical

PostgreSQL 9.4

· pg_hba.conf

· postgresql.conf wal_max_sender
0 10

· postgresql.conf wal_level logical

· postgresql.conf max_replication_slots
0 10

```
postgresql.conf
shared_preload_libraries=pglogical
shared_preload_libraries=pglogical
```

```
shared_preload_libraries = 'pglogical'
```

6 pglogical Extension

pglogical Extension
PostgreSQL 9.4 pglogical_origin
Extension

```
mydb=# SELECT version();
          version
-----
PostgreSQL 9.4.15 on x86_64-unknown-linux-gnu, compiled by gcc
(GCC) 4.8.5 20150623 (Red Hat 4.8.5-16), 64-bit
(1 row)
mydb=# CREATE EXTENSION pglogical_origin;
CREATE EXTENSION
mydb=# CREATE EXTENSION pglogical;
CREATE EXTENSION
mydb=# \dx pglogical
          List of installed extensions
  Name      | Version | Schema  | Description
-----+-----+-----+-----
pglogical  | 2.1.0   | pglogical | PostgreSQL Logical Replication
(1 row)
```

pglogical Extension

```
mydb=# SELECT version();
          version
-----
PostgreSQL 10.1 on x86_64-pc-linux-gnu, compiled by gcc (GCC)
4.8.5 20150623 (Red Hat 4.8.5-16), 64-bit
(1 row)
```

```
mydb=# CREATE EXTENSION pglogical;
CREATE EXTENSION
mydb=# \dx pglogical
```

```
              List of installed extensions
   Name      | Version | Schema  | Description
-----+-----+-----+-----
 pglogical   | 2.1.0   | pglogical | PostgreSQL Logical Replication
(1 row)
```

7 provider node

```
mydb=# SELECT pglogical.create_node(node_name := 'pg94provider', dsn
:= 'host=127.0.0.1 port=1921 dbname=mydb');
      create_node
-----
 3412564209
(1 row)
```

pglogical pglogical_node_info provider node

```
mydb=# SELECT * FROM pglogical.pglogical_node_info();
 node_id | node_name | sysid | dbname |
-----+-----+-----+-----+
 3412564209 | pg94provider | 6488108006998879651 | mydb |
 ".\x0Bg0\x02",",vg\x18\x02","Lg\x02"
(1 row)
```

8 provider node

```
mydb=# SELECT
pglogical.create_replication_set('insert_update_delete_notruncate',TR
UE,TRUE,TRUE,FALSE);
      create_replication_set
-----
 798613796
(1 row)
```



```

mydb=# SELECT * FROM pglogical.replication_set WHERE set_name =
'insert_update_delete_notruncate';
      set_id      | set_nodeid |      set_name      |
replicate_insert | replicate_update | replicate_delete |
replicate_truncate
-----+-----+-----+-----+
--
      798613796 | 3412564209 | insert_update_delete_notruncate | t
| t              | t              | f              |
(1 row)

```

9 provider node

```

mydb=# SELECT
pglogical.replication_set_add_all_tables('insert_update_delete_notrun
cate', ARRAY['public']);
      replication_set_add_all_tables
-----
t
(1 row)

```

10 PostgreSQL 10 subscription node

```

mydb=# SELECT pglogical.create_node(node_name := 'pg10subscriber',dsn
:= 'host= 127.0.0.1 port=1922 dbname=mydb');
      create_node
-----
      3105221948
(1 row)
mydb=# SELECT pglogical.create_subscription(subscription_name :=
'pg10subscription', provider_dsn := 'host=127.0.0.1 port=1921
dbname=mydb');
      create_subscription
-----
      1639099831
(1 row)

```

11 subscription node

```
mydb=# SELECT
pglogical.alter_subscription_synchronize('pg10subscription',FALSE);
alter_subscription_synchronize
-----
t
(1 row)
```

```

    "finished sync of table xxx for
subscriber pg10subscription"

```

```
mydb=# SELECT pglogical.show_subscription_status('pgl0subscription');
show_subscription_status
-----
      (pgl0subscription,replicating,pg94provider,"host=127.0.0.1
port=1921 dbname=mydb",pgl_mydb_pg94provider_pgl0subscription,
{insert_update_delete_nottruncate},{all})
      (1 row)
```

□□□□□□□□□□□□□□□□□□□□

12□□□□□□□□□□□□□□□□

```
mydb=# ALTER SYSTEM SET default_transaction_read_only = 'on';
ALTER SYSTEM
```

13

[illegible]

githubpglogical

[illegible]

15.4 备份与恢复

PostgreSQL 提供了多种备份与恢复工具。pg_dumpall 用于备份整个数据库系统，pglogical 用于逻辑复制，pg_upgrade 用于升级数据库系统。PostgreSQL 10 提供了 pg_upgrade 工具，用于升级数据库系统。PostgreSQL 10 提供了 pg_upgrade 工具，用于升级数据库系统。

PostgreSQL 10 提供了 pg_upgrade 工具，用于升级数据库系统。PostgreSQL 10 提供了 pg_upgrade 工具，用于升级数据库系统。

16 数据库

PostgreSQL数据库安装与配置
PostgreSQL数据库安装与配置

- PostgreSQL数据库world数据库安装
- GitHub数据库安装

数据库安装50个数据库安装
数据库安装PostgreSQL数据库安装
数据库安装PostgreSQL 9.0数据库安装

数据库安装
数据库安装

16.1 CREATE EXTENSION

```
make postgresql-world
$PGHOMESHARE/extension/
PostgreSQL
```

```
# make world
# make install-world
```

```

$PGHOMESHARE/extension/
```

```
$ ll $PGHOMESHARE/extension/
...
-rw-r--r-- 1 root root  794 Oct  6 10:20 pg_buffercache-
1.2.sql
-rw-r--r-- 1 root root  157 Oct  6 10:20
pg_buffercache.control
...
```

```
pg_buffercache.control
make world
```

```
make world PostgreSQL
CREATE
EXTENSION
```

```
CREATE EXTENSION [ IF NOT EXISTS ] extension_name
```

PostgreSQL 9.1

```
CREATE EXTENSION
```

SQL

```
pg_buffercache
```

```
postgres=# CREATE EXTENSION pg_buffercache;
CREATE EXTENSION
```

```
\dx
```

```
postgres=# \dx
List of installed extensions
Name | Version | Schema | Description
-----+-----+-----+-----
pg_buffercache | 1.3 | public | examine the shared buffer cache
plpgsql | 1.0 | pg_catalog | PL/pgSQL procedural language
(2 rows)
```

```
pg_extension
```

```

postgres=# SELECT * FROM pg_extension ;
extname      | extowner | extnamespace | extrelocatable |
extversion   | extconfig | extcondition
-----+-----+-----+-----+-----
plpgsql      |          | 10           | 11             | f
1.0          |          |              |                |
pg_buffercache |          | 10           | 2200           | t
1.3          |          |              |                |
(2 rows)

```

pg_available_extensions

```

postgres=# SELECT * FROM pg_available_extensions;
name          | default_version | installed_version |
comment
-----+-----+-----+-----
...
pg_buffercache | 1.3             | 1.3              |
examine the shared buffer cache
...
(43 rows)

```

pg_available_extensions

- name
- default_version
- installed_version

·comment□□□□□□□□

□□□□□□□□□□□□□□

16.2 pg_stat_statements

pg_stat_statements is a PostgreSQL extension that tracks the execution of SQL statements in a database. It provides a way to monitor the performance of SQL statements and identify slow or inefficient queries. The extension is installed by default in PostgreSQL 10 and later versions.

pg_stat_statements is a PostgreSQL extension that tracks the execution of SQL statements in a database. It provides a way to monitor the performance of SQL statements and identify slow or inefficient queries. The extension is installed by default in PostgreSQL 10 and later versions.

```
shared_preload_libraries = 'pg_stat_statements'      #
(change requires restart)
pg_stat_statements.max = 10000
pg_stat_statements.track = all
pg_stat_statements.track_utility = on
pg_stat_statements.save = on
```

shared_preload_libraries is a PostgreSQL configuration parameter that specifies the shared libraries to be loaded at database startup. It is used to load extensions that are not installed by default. The parameter is set in the postgresql.conf file.

pg_stat_statements is a PostgreSQL extension that tracks the execution of SQL statements in a database. It provides a way to monitor the performance of SQL statements and identify slow or inefficient queries. The extension is installed by default in PostgreSQL 10 and later versions.

·pg_stat_statements.max
pg_stat_statements SQL
5000 SQL

·pg_stat_statements.track
pg_stat_statements SQL top
SQL all SQL all

·pg_stat_statements.track_utility
pg_stat_statements SELECT
UPDATE DELETE INSERT SQL on

·pg_stat_statements.save
pg_stat_statements SQL
off SQL
on

```
[postgres@pghost1 pg_root]$ pg_ctl restart -m fast
```

postgres postgres
pg_stat_statements

```
[postgres@pghost1 loadtest]$ psql postgres
psql (10.0)
Type "help" for help.
```

```
postgres=# CREATE EXTENSION pg_stat_statements ;
CREATE EXTENSION
```

pg_stat_statements
pg_stat_statements_reset
pg_stat_statements 16-1

16-1 pg_stat_statements

名 称	类 型	关 联 信 息	描 述
userid	oid	pg_authid.oid	执行此 SQL 的用户 OID
dbid	oid	pg_database.oid	此 SQL 执行的数据库 OID
queryid	bigint		此 SQL 的编号
query	text		此 SQL 内容
calls	bigint		此 SQL 调用次数
total_time	double precision		此 SQL 执行的总时间，单位毫秒
min_time	double precision		此 SQL 执行的最小时间，单位毫秒
max_time	double precision		此 SQL 执行的最大时间，单位毫秒
mean_time	double precision		此 SQL 执行的平均时间，单位毫秒
rows	bigint		此 SQL 影响的数据行
shared_blks_hit	bigint		此 SQL 命中的共享内存数据块数
shared_blks_read	bigint		此 SQL 读取的共享内存数据块数
shared_blks_dirtied	bigint		此 SQL 产生的共享内存数据脏块数量
shared_blks_written	bigint		此 SQL 写入的共享内存数据块数

pg_stat_statements
pgbench SQL
SQL pg_stat_statements

tran_per1.sql

```
\set v_id random(1,10000000)

SELECT name FROM test_per1 WHERE id=:v_id;
UPDATE test_per2 SET flag='1' WHERE id=:v_id;
```

pgbench

```
$ pgbench -c 4 -T 120 -d postgres -U postgres -n N -M
prepared -f tran_per1.sql >
tran_per1_4.out 2>&1 &
```

pg_stat_statements

SQL SQL SQL calls

```
postgres=# SELECT userid,dbid,queryid,query,calls,
      total_time,min_time,max_time,mean_time,rows
      FROM pg_stat_statements ORDER BY calls DESC LIMIT 2;
-[ RECORD 1 ]-----
userid      | 10
dbid        | 13158
queryid     | 1758110241
query       | UPDATE test_per2 SET flag=$2 WHERE id=$1
calls       | 812691
total_time  | 83537.8480209997
min_time    | 0.015113
max_time    | 6583.11359
mean_time   | 0.102791649004361
rows        | 812691
-[ RECORD 2 ]-----
userid      | 10
dbid        | 13158
```

queryid		2354798721
query		SELECT name FROM test_per1 WHERE id=\$1
calls		812691
total_time		28687.6311289994
min_time		0.008581
max_time		76.785687
mean_time		0.0352995555863171
rows		812691

SQL

```
postgres=# SELECT userid,dbid,queryid,query,calls,
total_time,min_time,max_time,mean_time,rows
FROM pg_stat_statements ORDER BY mean_time DESC LIMIT 1;
-[ RECORD1 ]-----
userid      | 10
dbid        | 13158
queryid     | 125206108
query       | select pg_sleep($1)
calls       | 1
total_time  | 6019.78287
min_time    | 6019.78287
max_time    | 6019.78287
mean_time   | 6019.78287
```

SQLmean_time
SQL

pg_stat_statements
pg_stat_statements_reset
SQL

```
postgres=# SELECT COUNT(*) FROM pg_stat_statements ;
count
```

```
-----
      20
(1 row)
```

```
postgres=# SELECT pg_stat_statements_reset();
pg_stat_statements_reset
```

```
-----
(1 row)
```

```
postgres=# SELECT COUNT(*) FROM pg_stat_statements ;
count
```

```
-----
      2
(1 row)
```

```
pg_stat_statements
CREATE
EXTENSION
```

16.3 auto_explain

EXPLAIN EXPLAIN ANALYZE
SQL SQL SQL SQL SQL
SQL SQL SQL SQL SQL
SQL PostgreSQL
SQL

```

    auto_explain SQL
    SQL
    SQL

```

```
auto_explain=on pg_stat_statements=on
postgresql.conf postgresql.conf
shared_preload_libraries=auto_explain
```

```
shared_preload_libraries = 'pg_stat_statements,auto_explain'
auto_explain.log_min_duration = 0
auto_explain.log_analyze = on
auto_explain.log_buffers = off
```

```

    auto_explain
postgresql.conf

```

```
·auto_explain.log_min_duration=SQL
SQL
SQL
0SQL-1
```

100ms 100 SQL
100

·auto_explain.log_analyze
SQL ANALYZE EXPLAIN
ANALYZE off

·auto_explain.log_buffers
SQL EXPLAIN
BUFFERS off

auto_explain
EXPLAIN EXPLAIN

```
[postgres@pgghost1 pg_log]$ pg_ctl restart -m fast
```

SQL
SQL

```
postgres=# SELECT * FROM test_per1 WHERE id=1;
 id | name | create_time
-----+-----
  1 | 1_per1 | 2017-09-04 21:21:44
(1 row)
```

```
[postgres@pgghost1 pg_log]$ tail -f postgresql-2017-10-22_201048.csv
2017-10-22 20:18:25.927 CST,"postgres","postgres",24856,"[local]",59ec8c8d.6118,1,"SELECT",2017-10-22 20:18:21 CST,5/367,0,LOG,00000,"duration: 0.049 ms  plan: Query Text: SELECT * FROM test_perl WHERE id=1; Index Scan using test_perl_pkey on test_perl (cost=0.43..4.45 rows=1 width=24) (actual time=0.023..0.023 rows=1 loops=1) Index Cond: (id = 1)",,,,,,,,,,"psql"
```

□□□□□□□□□□SQL□□□□SQL□□□□□□□□

```
auto_explain.log_min_duration 0
auto_explain.log_max_duration 100
auto_explain.log_buffer_size 1024
auto_explain.log_sql
```

auto_explain PostgreSQL
SQL
SQL

16.4 pg_prewarm

pg_prewarm is a PostgreSQL 9.4 extension that allows you to prewarm database buffers. It is installed by default in PostgreSQL 9.4 and later versions. The extension is used to prewarm buffers for a specific database or schema, which can improve performance by reducing the time it takes to load data from disk into memory.

postgres=# CREATE EXTENSION pg_prewarm;
CREATE EXTENSION

```
[postgres@pghost1 ~]$ psql mydb postgres
psql (10.0)
Type "help" for help.
```

```
mydb=# CREATE EXTENSION pg_prewarm;
CREATE EXTENSION
```

pg_prewarm is a PostgreSQL 9.4 extension that allows you to prewarm database buffers. It is installed by default in PostgreSQL 9.4 and later versions. The extension is used to prewarm buffers for a specific database or schema, which can improve performance by reducing the time it takes to load data from disk into memory.

```
pg_prewarm(regclass, mode text default 'buffer', fork text
default 'main',
    first_block int8 default null,
    last_block int8 default null) RETURNS int8
```

pg_prewarm is a PostgreSQL 9.4 extension that allows you to prewarm database buffers. It is installed by default in PostgreSQL 9.4 and later versions. The extension is used to prewarm buffers for a specific database or schema, which can improve performance by reducing the time it takes to load data from disk into memory.

```
·regclass
```

```
·mode[]prefetch[]
[]read[]
[]buffer[]
```

```
·fork[ ][ ][ ][ ][ ]main[ ][ ][ ][ ][ ][ ][ ][ ][ ][ ]
```

```

    ·first_block[0] = null
    [0] = 0

```

```
·last_block[0][0][0][0][0][0][0][0]null[0][0][0][0]
[0][0][0][0][0][0][0][0]
```

```
mydbt_pre200

```

```
[postgres@pghost1 ~]$ psql mydb pguser
mydb=> CREATE TABLE t_pre(id int4 PRIMARY KEY,
        info text,
        create_time timestamp(0) without time zone);
CREATE TABLE

mydb=> INSERT INTO t_pre(id,info,create_time)
        SELECT n,n||'_pre',clock_timestamp() FROM
generate_series(1,2000000) n ;
INSERT 0 2000000
```

It pre

```
mydb=> SELECT pg_prewarm('t_pre','buffer');
pg_prewarm
-----
12739
(1 row)
```

12739

t_pre pg_class

```
mydb=> SELECT relname,relpages FROM pg_class WHERE
relname='t_pre';
relname | relpages
-----+-----
t_pre   |      12739
(1 row)
```

relpages ANALYZE


```
mydb=> SELECT pg_prewarm('t_pre','prefetch');
pg_prewarm
-----
12739
(1 row)
```

```
mydb=> SELECT pg_prewarm('t_pre','read');
pg_prewarm
```

12739
(1 row)

pg_prewarm
pg_prewarm
pg_prewarm

pgfincore
PostgreSQL9.4 pg_prewarm
pgfincore
<https://github.com/klando/pgfincore>

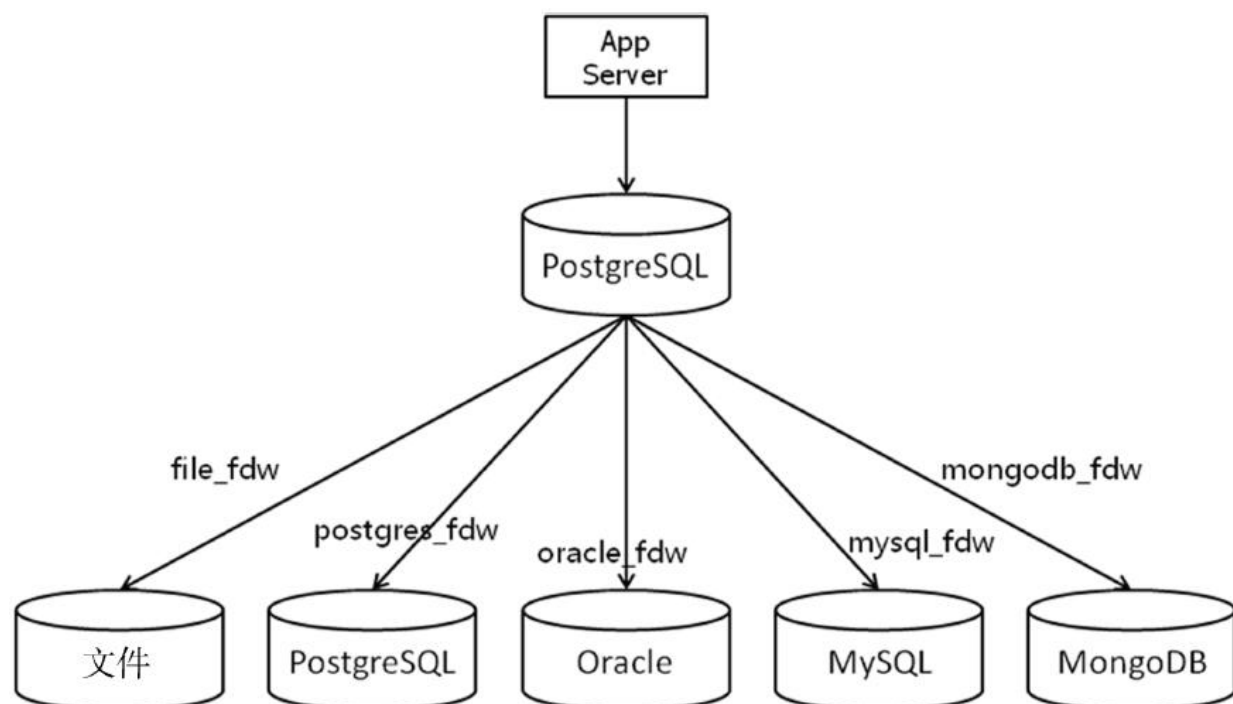
 pg_prewarm
pg_prewarm

16.5 file_fdw

file_fdw PostgreSQL 外部数据源
file_fdw SQL/MED SQL
Management of External Data

16.5.1 SQL/MED

SQL/MED PostgreSQL 外部数据源
PostgreSQL SQL 外部数据源
Oracle dblink SQL/MED
16-1



16-1 SQL/MED

データベースと外部データ

- ・PostgreSQLデータベースと外部データ
外部データ形式: csv, text

- ・PostgreSQLデータベースと外部データ
PostgreSQL, Oracle, MySQL, SQL Server

- ・MongoDB, Redis, Cassandra
NoSQLデータベース

- ・Elastic Search, Hadoop

PostgreSQLデータベースと外部データ
外部データ形式: csv, text, binary

16.5.2 file_fdw

PostgreSQLデータベースと外部データ
file_fdwデータベースと外部データ
text, csv, binary

file_fdwデータベースと外部データ

1 安装file_fdw扩展

2 安装foreign server扩展并创建foreign server
扩展

3 安装并配置file_fdw扩展

4 测试

在postgres数据库中安装file_fdw扩展
mydb=# file_fdw

```
[postgres@pghost1 ~]$ psql mydb postgres
psql (10.0)
Type "help" for help.
```

```
mydb=# CREATE EXTENSION file_fdw;
CREATE EXTENSION
```

安装file_fdw扩展并创建fs_file扩展

```
mydb=# CREATE SERVER fs_file FOREIGN DATA WRAPPER file_fdw;
CREATE SERVER
```

在fs_file扩展中安装des扩展
在des扩展中安装des扩展
在des扩展中安装des扩展

```
mydb=# \des
```

List of foreign servers

Name	Owner	Foreign-data wrapper
------	-------	----------------------

-----+-----+-----

```
fs_file | postgres | file_fdw
```

(1 row)

IP

```

/home/postgres/script/file1.txt
tab

```

1	a
2	b

```

file_fdw

```

```
mydb=# CREATE FOREIGN TABLE ft_file1(
        id int4,
        flag text
    ) SERVER fs_file
      OPTIONS (filename
'/home/postgres/script/file1.txt',format 'text');
CREATE FOREIGN TABLE
```

```
OPTIONS file_fdw fdw
file_fdw
```

·filename[문자열]
[문자열]

·format[문자열]text|csv|
binary[문자열]text|

·header[문자열]csv[문자열]
[문자열]

·delimiter[문자열]text[문자열]
[문자열]tab[문자열]

·encoding[문자열]

[문자열]filename[문자열]copy[문자열]

[문자열]
[문자열]

```
mydb=# SELECT * FROM ft_file1 ;
      id | flag
-----+-----
       1 | a
       2 | b
(2 rows)
```

[문자열]
[문자열]file_fdw[문자열]

PostgreSQL

file_fdw
INSERT/UPDATE/DELETE

\det

```
mydb=# \det
      List of foreign tables
 Schema | Table   | Server
-----+-----+-----
 public | ft_file1 | fs_file
(1 row)
```

16.5.3 file_fdw

file_fdw PostgreSQL
csv
postgresql.conf

```
log_destination = 'csvlog'
logging_collector = on
log_directory = 'pg_log'
log_filename = 'postgresql-%Y-%m-%d_%H%M%S.log'
```

·log_destination postgresql stderr csvlog syslog csv logging_collector on

·logging_collector on PostgreSQL stderr PostgreSQL

·log_directory logging_collector PostgreSQL \$PGDATA

·log_filename logging_collector PostgreSQL

\$PGDATA/pg_log/postgresql-2017-10-26_000000.csv

```
...
2017-10-26 14:29:16.298 CST,"postgres","mydb",18701,"
[local]",59f1762a.490d,31,"DELETE",2017-10-26 13:44:10
CST,5/16770,0,ERROR,0A000,"cannot delete from foreign table
""ft_file1""",,,,,,"DELETE FROM ft_file1 ;",,,,"psql"
...
```

PostgreSQL

```

CREATE FOREIGN TABLE ft_pglog (
    log_time timestamp(0) without time zone,
    user_name text,
    database_name text,
    process_id integer,
    connection_from text,
    session_id text,
    session_line_num bigint,
    command_tag text,
    session_start_time timestamp with time zone,
    virtual_transaction_id text,
    transaction_id bigint,
    error_severity text,
    sql_state_code text,
    message text,
    detail text,
    hint text,
    internal_query text,
    internal_query_pos integer,
    context text,
    query text,
    query_pos integer,
    location text,
    application_name text
) SERVER fs_file
OPTIONS ( filename
'/database/pg10/pg_root/pg_log/postgresql-2017-10-
26_000000.csv', format 'csv' );

```

```

OPTIONS
csv

```

```

mydb=# SELECT * FROM ft_pglog WHERE error_severity='ERROR'
AND command_tag='DELETE';

```

```

-[ RECORD 1 ]-----+-----

```

```

log_time          | 2017-10-26 14:29:16
user_name         | postgres
database_name     | mydb
process_id        | 18701

```

connection_from		[local]
session_id		59f1762a.490d
session_line_num		31
command_tag		DELETE
session_start_time		2017-10-27 03:44:10+08
virtual_transaction_id		5/16770
transaction_id		0
error_severity		ERROR
sql_state_code		0A000
message		cannot delete from foreign table
"ft_file1"		
detail		
hint		
internal_query		
internal_query_pos		
context		
query		DELETE FROM ft_file1 ;
query_pos		
location		
application_name		psql

postgresql-2017-10-
 26_000000.csv
 ft_pglog

16.6 postgres_fdw

通过file_fdw将PostgreSQL数据库以csv格式导入到postgres_fdw数据库，并连接到PostgreSQL数据库的Oracle dblink。

16.6.1 postgres_fdw

通过postgres_fdw将16-2

16-2 postgres_fdw

角 色	主机名	IP	端 口	库 名	用户名	版 本
本地库	pghost1	192.168.28.74	1921	mydb	pguser	PostgreSQL10
远程库	pghost3	192.168.28.76	1923	des	pguser	PostgreSQL10

通过pghost1 mydb postgres_fdw
通过pghost1 mydb
pghost3 des

postgres_fdw

1 postgres_fdw

2 foreign server
通过


```
'pghost3', port '1923', dbname 'des');  
CREATE SERVER
```

```
fs_postgres_pghost3  
OPTIONS PostgreSQL  
PostgreSQL
```

```
mydb=> CREATE USER MAPPING FOR pguser  
        SERVER fs_postgres_pghost3 OPTIONS (user 'pguser',  
        password 'pguser');  
CREATE USER MAPPING
```

```
FOR OPTIONS  
PostgreSQL  
PostgreSQL IP  
PostgreSQL
```

```
pghost4 des  

```

```
[des@pghost3 ~]$ psql des pguser  
des=> CREATE TABLE t_fdw1 (id int4,info text);  
CREATE TABLE  
  
des=> INSERT INTO t_fdw1 (id, info ) VALUES (1,'a'),  
('2','b');  
INSERT 0 2
```

pgghost1

```
mydb=> CREATE FOREIGN TABLE ft_t_fdw1 (  
    id      int4,  
    info   text  
) SERVER fs_postgres_pghost3 OPTIONS (schema_name 'pguser',  
table_name 't_fdw1');
```

OPTIONS schema_name
table_name

postgres_fdw
ft_t_fdw1

```
mydb=> SELECT * FROM ft_t_fdw1 ;  
ERROR:  could not connect to server "fs_postgres_pghost3"  
DETAIL:  FATAL:  no pg_hba.conf entry for host  
"20.26.28.74", user "pguser", database "des"
```

pg_hba.conf
pghost3 \$PGDATA/pg_hba.conf

host	des	pguser	20.26.28.0/24	md5
------	-----	--------	---------------	-----

pg_ctl reload

pghost1 ft_t_fdw1

```
mydb=> SELECT * FROM ft_t_fdw1 ;
   id | info
-----+-----
    1 | a
    2 | b
(2 rows)
```

pgghost1 mydb
pgghost3 des

16.6.2 postgres_fdw

postgres_fdw
PostgreSQL9.3 postgres_fdw

pgghost1

```
mydb=> INSERT INTO ft_t_fdw1 (id,info) VALUES (3,'c');
INSERT 0 1
mydb=> SELECT * FROM ft_t_fdw1 WHERE id=3;
   id | info
-----+-----
    3 | c
(1 row)
```

pgghost3 des

```
des=> SELECT * FROM t_fdw1 WHERE ID=3;
   id | info
```

```
-----+-----
      3 | c
(1 row)
```

pgghost1 des

pgghost1

```
mydb=> UPDATE ft_t_fdw1 SET info='ccc' WHERE ID=3;
UPDATE 1
mydb=> SELECT * FROM ft_t_fdw1 WHERE id=3;
      id | info
-----+-----
      3 | ccc
(1 row)
```

pgghost3 des

```
des=> SELECT * FROM t_fdw1 WHERE ID=3;
      id | info
-----+-----
      3 | ccc
(1 row)
```

DELETE pgghost1

DELETE

```
mydb=> DELETE FROM ft_t_fdw1 WHERE id=3;
DELETE 1
```

pgghost3des

```
des=> SELECT * FROM t_fdw1 WHERE ID=3;
      id | info
-----+-----
(0 rows)
```



postgres_fdw
PostgreSQL9.3

16.6.3 postgres_fdw

PostgreSQL10postgres_fdw
PostgreSQL
10postgres_fdw10

PostgreSQL10

pgghost3 des

```
[des@pgghost3 ~]$ psql des pguser
psql (10.0)
Type "help" for help.
```

```
des=> CREATE TABLE t_fdw2(id int4,flag int4);
CREATE TABLE
```

```
des=> INSERT INTO t_fdw2(id,flag) SELECT n,mod(n,3) FROM
generate_series(1,100000) n;
INSERT 0 100000
```

pgghost1 mydb

```
mydb=> CREATE FOREIGN TABLE ft_t_fdw2 (
    id    int4,
    flag  int4
) SERVER fs_postgres_pgghost3 OPTIONS (schema_name 'pguser',
table_name 't_fdw2');
```

SQL

```
mydb=> SELECT flag,count(*) FROM ft_t_fdw2 GROUP BY flag
ORDER BY flag;
   flag | count
-----+-----
       0 | 33333
       1 | 33334
       2 | 33333
(3 rows)
```

Query Plan

```
mydb=> EXPLAIN (ANALYZE on,VERBOSE on)
        SELECT flag,count(*) FROM ft_t_fdw2 GROUP BY flag ORDER
        BY flag;
QUERY PLAN
-----
          Sort  (cost=167.52..168.02 rows=200 width=12) (actual
time=18.888..18.888 rows=3 loops=1)
            Output: flag, (count(*))
            Sort Key: ft_t_fdw2.flag
            Sort Method: quicksort  Memory: 25kB
            -> Foreign Scan  (cost=114.62..159.88 rows=200
width=12) (actual time= 18.877..18.878 rows=3 loops=1)
              Output: flag, (count(*))
              Relations: Aggregate on (pguser.ft_t_fdw2)
              Remote SQL: SELECT flag, count(*) FROM
pguser.t_fdw2 GROUP BY flag
              Planning time: 0.109 ms
              Execution time: 19.496 ms
(10 rows)
```

Query Plan Details

·Remote SQL: SELECT flag, count(*) FROM pguser.t_fdw2 GROUP BY flag ORDER BY flag

·Relation Aggregate: pguser.ft_t_fdw2

·Foreign Scan: "rows=3" Foreign Scan: pguser.t_fdw2

·Sort

SQL19.496ms
SQL

Remote SQL: SELECT flag, count(*) FROM pguser.t_fdw2 GROUP BY flag

PostgreSQL9.6

9.616-3

16-3 postgres_fdw

角 色	主机名	IP	端 口	库 名	用户名	版 本
本地库	pghost1	192.168.28.74	1922	mydb	pguser	PostgreSQL9.6
远程库	pghost3	192.168.28.76	1924	des	pguser	PostgreSQL9.6

PostgreSQL9.6postgres_fdw
16.6.1

```
mydb=> EXPLAIN (ANALYZE on,VERBOSE on)
        SELECT flag,count(*) FROM ft_t_fdw2 GROUP BY flag
ORDER BY flag;
QUERY PLAN
```

```

Sort (cost=324.43..324.93 rows=200 width=12) (actual
time=354.879..354.880 rows=3 loops=1)
  Output: flag, (count(*))
  Sort Key: ft_t_fdw2.flag
  Sort Method: quicksort Memory: 25kB
  -> HashAggregate (cost=314.78..316.78 rows=200
width=12) (actual time= 354.863..354.863 rows=3 loops=1)
    Output: flag, count(*)
    Group Key: ft_t_fdw2.flag
    -> Foreign Scan on pguser.ft_t_fdw2
(cost=100.00..285.53 rows=5851 width=4) (actual
time=1.014..333.186 rows=100000 loops=1)
      Output: id, flag
      Remote SQL: SELECT flag FROM pguser.t_fdw2
    Planning time: 0.152 ms
    Execution time: 355.649 ms
(12 rows)

```

355 PostgreSQL10
 17 Remote SQL SQL

Remote SQL: SELECT flag FROM pguser.t_fdw2

Foreign Scan
 rows=100000
 pghost1 pghost1

PostgreSQL10 SQL Remote
 SQL SQL

Remote SQL: SELECT flag, count(*) FROM pguser.t_fdw2 GROUP BY flag

PostgreSQL10
postgres_fdw

16.7 Citus

Citus is a distributed database system for PostgreSQL. It is designed for B2B applications that require high performance and scalability. Citus is built on top of PostgreSQL and uses the pg_shard extension to manage the distributed data. Citus is a good choice for applications that require high performance and scalability.

16.7.1 Citus

1. Introduction

Citus is a distributed database system for PostgreSQL. It is designed for B2B applications that require high performance and scalability. Citus is built on top of PostgreSQL and uses the pg_shard extension to manage the distributed data. Citus is a good choice for applications that require high performance and scalability.

2. Features

Citus is a distributed database system for PostgreSQL. It is designed for B2B applications that require high performance and scalability. Citus is built on top of PostgreSQL and uses the pg_shard extension to manage the distributed data. Citus is a good choice for applications that require high performance and scalability.

3. Citus

Citus is a distributed database for PostgreSQL. It is designed for B2B applications that require high performance and scalability. Citus allows you to run PostgreSQL on a cluster of machines, distributing data and queries across the nodes. This makes it ideal for large-scale data processing and analytics.

4. Citus

Citus is a distributed database for PostgreSQL. It is designed for B2B applications that require high performance and scalability. Citus allows you to run PostgreSQL on a cluster of machines, distributing data and queries across the nodes. This makes it ideal for large-scale data processing and analytics.

16.7.2 Citus

`yum`

```
# yum install citus
curl https://install.citusdata.com/community/rpm.sh | sudo
bash
sudo yum install -y citus72_10
```

`pgsql`

```
sudo su - postgres
export PATH=$PATH:/usr/pgsql-10/bin
```

mkdir -p /pgdata/citus_cluster/coordinator
/pgdata/citus_cluster/worker1 /pgdata/citus_cluster/worker2

pg_ctl initdb -D /pgdata/citus_cluster/coordinator/

```
/usr/pgsql-10/bin/initdb -D /pgdata/citus_cluster/coordinator/  
/usr/pgsql-10/bin/initdb -D /pgdata/citus_cluster/worker1/  
/usr/pgsql-10/bin/initdb -D /pgdata/citus_cluster/worker2/
```

echo "shared_preload_libraries = 'citus'" >> /pgdata/citus_cluster/coordinator/postgresql.conf

```
-- shared_preload_libraries = 'citus'  
echo "shared_preload_libraries = 'citus'" >>  
/pgdata/citus_cluster/coordinator/postgresql.conf  
echo "shared_preload_libraries = 'citus'" >>  
/pgdata/citus_cluster/worker1/postgresql.conf  
echo "shared_preload_libraries = 'citus'" >>  
/pgdata/citus_cluster/worker2/postgresql.conf
```

FATAL: Citus has to be loaded first
HINT: Place citus at the beginning of

shared_preload_libraries.

□□□□□□□□□□

```
/usr/pgsql-10/bin/pg_ctl -D
/pgdata/citus_cluster/coordinator -o "-p 9700" -l
coordinator_logfile start
/usr/pgsql-10/bin/pg_ctl -D /pgdata/citus_cluster/worker1 -o
"-p 9701" -l worker1_logfile start
/usr/pgsql-10/bin/pg_ctl -D /pgdata/citus_cluster/worker2 -o
"-p 9702" -l worker2_logfile start
```

□Coordinator□□□□□□□□□□

```
/usr/pgsql-10/bin/psql -p 9700 postgres -c "CREATE DATABASE
mydb;"
NOTICE: Citus partially supports CREATE DATABASE for
distributed databases
DETAIL: Citus does not propagate CREATE DATABASE command to
workers
HINT: You can manually create a database and its extensions
on workers.
CREATE DATABASE
```

□Coordinator□□□□□□□□□□
Coordinator□□□□□□□□□□□□□□Worker□□□□
□□□□□□□□□□□□□□Worker□□□□□□□□□□□□□□□□
□□□□□Citus Extension□

□Worker□□□□□□□□□□□□□□

```
/usr/pgsql-10/bin/psql -p 9701 postgres -c "CREATE DATABASE mydb;"
/usr/pgsql-10/bin/psql -p 9702 postgres -c "CREATE DATABASE mydb;"
```

安装 Citus Extension

```
/usr/pgsql-10/bin/psql -p 9700 mydb -c "CREATE EXTENSION citus;"
/usr/pgsql-10/bin/psql -p 9701 mydb -c "CREATE EXTENSION citus;"
/usr/pgsql-10/bin/psql -p 9702 mydb -c "CREATE EXTENSION citus;"
```

16.7.3 Citus

Citus Work Work
Work Work

Citus Worker

```
/usr/pgsql-10/bin/psql -p 9700 mydb -c "SELECT * FROM master_add_node('127. 0.0.1', 9701);"
 nodeid | groupid | nodename | nodeport | noderack | hasmetadata | isactive
-----+-----+-----+-----+-----+-----+-----
1 | 1 | 127.0.0.1 | 9701 | default | f
(1 row)
/usr/pgsql-10/bin/psql -p 9700 mydb -c "SELECT * FROM master_add_node('127. 0.0.1', 9702);"
 nodeid | groupid | nodename | nodeport | noderack |
```

```
ERROR: failed to assign 2 task(s) to worker nodes
```

```
CoordinatorWorkerWorkerWorkerWorker
```

```
SELECT master_activate_node('127.0.0.1', 9702);
```

```
CoordinatorWorkerWorkerWorker
```

```
/usr/pgsql-10/bin/psql -p 9700 mydb -c " SELECT * FROM
master_get_active_worker_nodes();"
 node_name | node_port
-----+-----
 127.0.0.1 |      9701
 127.0.0.1 |      9702
(2 rows)
```

16.7.4 配置

```
CoordinatorWorkerWorkerWorker
```

```
CREATE TABLE table_name (
  user_id integer NOT NULL,
  name character varying(32)
);
```

```
CoordinatorWorkerWorkerWorker
```

```
mydb=# SELECT create_distributed_table('table_name',
'user_id');
      create_distributed_table
-----
(1 row)
```

Citus 32
 Citus
 CPU x x
 Citus

```
set citus.shard_count = 64;
```

```

    Worker.join()
    Coordinator.create_reference_table()
    Worker

```

```
CREATE TABLE rt(id serial primary key,ival int);
SELECT create_reference_table('rt');
```

```

1 upgrade_to_reference_table API

```

16.7.5 Citus

·citus.shard_replication_factor
integer

1
2

·citus.shard_count integer

32
Citush
Worker
3
Worker1
2
Worker2
1

·citus.task_executor_typeenum

real-time
task-tracker
real-time
task-tracker
worker

Citus
postgresql.conf
idle_in_transaction_session_timeout

```
SELECT get_shard_id_for_distribution_column('table_name',
DK_value);
```

create_distributed_table

```
mydb=# \df+ create_distributed_table
```

```
List of functions
```

```
-[ RECORD 1 ]-----+-----
Schema                | pg_catalog
Name                  | create_distributed_table
Result data type      | void
Argument data types   | table_name regclass,
distribution_column text, distribution_type
citius.distribution_type DEFAULT
'hash'::citius.distribution_type, colocate_with text DEFAULT
'default'::text
Type                  | normal
Volatility             | volatile
Parallel              | unsafe
Owner                 | postgres
Security              | invoker
Access privileges     |
Language              | c
Source code           | create_distributed_table
Description           | creates a distributed table
```

```
mydb=# \dT+ citius.distribution_type
```

```
List of data types
```

```
-[ RECORD 1 ]-----+-----
Schema                | citius
Name                  | citius.distribution_type
Internal name         | distribution_type
Size                  | 4
```

Elements	hash	+
	range	+
	append	
Owner	postgres	
Access privileges		
Description		

Citus

 Citus

```
SET citus.explain_all_tasks = 'TRUE';
```

DDL

 Citus

 ALTER TABLE ADD

 COLUMN

 ALTER COLUMN

 DROP COLUMN

 ALTER COLUMN TYPE

 RENAME COLUMN

 PostgreSQL

```
ALTER TABLE tbl ADD COLUMN col INT;
NOTICE: using one-phase commit for distributed DDL commands
HINT: You can enable two-phase commit for extra safety with:
SET citus.multi_shard_commit_protocol TO '2pc'
```

Citus

 DDL

 DDL

```
SET citus.multi_shard_commit_protocol TO '2pc';
ALTER TABLE tbl ADD COLUMN col INT;
ALTER TABLE test ALTER COLUMN ival3 SET DEFAULT 1;
```

PostgreSQL Citus
distributed PostgreSQL

```
ERROR: creating index without a name on a distributed table  
is currently unsupported
```

PostgreSQL

```
CREATE INDEX CONCURRENTLY idx_name ON table_name  
(col1...col_n);
```

citus
distributed PostgreSQL

```
DROP INDEX umch_user_id_name;  
NOTICE: using one-phase commit for distributed DDL commands  
HINT: You can enable two-phase commit for extra safety with:  
SET citus.multi_shard_commit_protocol TO '2pc'
```

VACUUM ANALYZE

```
VACUUM table_name;  
ANALYZE table_name;
```

`cituser@cituser:~$ sudo -u postgres psql -c 'ALTER DATABASE citus SET VERBOSE;'`
`cituser@cituser:~$ sudo -u postgres psql -c 'VACUUM citus;'`

ERROR: the VERBOSE option is currently unsupported in distributed VACUUM commands

`cituser@cituser:~$ sudo -u postgres psql -c 'CREATE EXTENSION citus;'`
`cituser@cituser:~$ sudo -u postgres psql -c 'CREATE EXTENSION PostGIS;'`
`cituser@cituser:~$ sudo -u postgres psql -c 'CREATE EXTENSION hll;'`
`cituser@cituser:~$ sudo -u postgres psql -c 'CREATE EXTENSION jsonb;'`

`cituser@cituser:~$ sudo -u postgres psql -c 'ALTER SYSTEM SET shared_preload_libraries = citus;'`
`cituser@cituser:~$ sudo -u postgres psql -c 'RELOAD;'`

`cituser@cituser:~$ sudo -u postgres psql -c 'ALTER SYSTEM SET citus_coordinator_node = citus;'`
`cituser@cituser:~$ sudo -u postgres psql -c 'ALTER SYSTEM SET citus_worker_node = citus;'`

`cituser@cituser:~$ sudo -u postgres psql -c 'ALTER SYSTEM SET citus_max_parallel_workers_per_gather = 4;'`

```
yum --showduplicates list citus72_10
Available Packages
citus72_10.x86_64      7.2.0.citus-1.el6      citusdata_community
citus72_10.x86_64      7.2.1.citus-1.el6      citusdata_community
```

`cituser@cituser:~$ sudo -u postgres psql -c 'ALTER SYSTEM SET citus_max_parallel_workers_per_gather = 4;'`

```
yum install -y 7.2.1.citus-1.el6
```

如何安装PostgreSQL

如何安装yum

16.8 数据库

PostgreSQL数据库是一个开源的数据库系统，它支持SQL语言，并且具有强大的安全性和性能。PostgreSQL数据库可以用于各种应用程序，从简单的Web应用到复杂的企业级应用。

PostgreSQL数据库的安装和配置相对简单，可以通过以下步骤完成：1. 下载PostgreSQL数据库的源代码。2. 编译并安装PostgreSQL数据库。3. 配置PostgreSQL数据库的参数文件。4. 启动PostgreSQL数据库服务。5. 创建数据库用户和角色。6. 创建数据库和表。7. 导入数据。8. 使用GitHub仓库中的脚本进行自动化部署。

17 Oracle 与 PostgreSQL

SQL/MED PostgreSQL
PostgreSQL
PostgreSQL
Oracle PostgreSQL

Oracle PostgreSQL
PostgreSQL oracle_fdw
Oracle 30 4
30G 100GB
PostgreSQL 9.2 9.3

Oracle PostgreSQL
PostgreSQL

17.1 資料庫

資料庫是儲存和管理數據的系統，它允許用戶訪問和修改數據。資料庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。

- 資料庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。數據庫引擎負責處理數據的存取和修改，數據庫管理系統負責管理數據庫的運行和維護，數據庫用戶負責訪問和修改數據。

- 資料庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。數據庫引擎負責處理數據的存取和修改，數據庫管理系統負責管理數據庫的運行和維護，數據庫用戶負責訪問和修改數據。PostgreSQL、Oracle、MySQL 是常見的數據庫系統。

- 資料庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。數據庫引擎負責處理數據的存取和修改，數據庫管理系統負責管理數據庫的運行和維護，數據庫用戶負責訪問和修改數據。Oracle、PostgreSQL、MySQL 是常見的數據庫系統。Oracle 數據庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。

- 資料庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。數據庫引擎負責處理數據的存取和修改，數據庫管理系統負責管理數據庫的運行和維護，數據庫用戶負責訪問和修改數據。DBA 是數據庫管理員的縮寫。

- 資料庫系統通常由數據庫引擎、數據庫管理系統（DBMS）和數據庫用戶組成。數據庫引擎負責處理數據的存取和修改，數據庫管理系統負責管理數據庫的運行和維護，數據庫用戶負責訪問和修改數據。DBA 是數據庫管理員的縮寫。數據庫管理員負責管理數據庫的運行和維護，包括數據庫的安裝、配置、性能優化、備份和恢復等。

A diagram showing a sequence of 30 boxes arranged in three rows. The first row has 10 boxes, the second row has 10 boxes, and the third row has 10 boxes. A small square is placed to the left of the first box in the first row.

□□□□□□□□□□DBA□□□□□□□□□□

17.2 数据类型

Oracle与PostgreSQL数据类型对照表
本表列出了Oracle与PostgreSQL数据库中的数据类型对照关系。表中列出的数据类型在两个数据库中都是有效的。表中列出的数据类型在两个数据库中的名称可能不同，但它们的含义是相同的。表中列出的数据类型在两个数据库中的名称可能不同，但它们的含义是相同的。表中列出的数据类型在两个数据库中的名称可能不同，但它们的含义是相同的。

1. 数据类型

本表列出了Oracle与PostgreSQL数据库中的数据类型对照关系。表中列出的数据类型在两个数据库中都是有效的。表中列出的数据类型在两个数据库中的名称可能不同，但它们的含义是相同的。表中列出的数据类型在两个数据库中的名称可能不同，但它们的含义是相同的。表中列出的数据类型在两个数据库中的名称可能不同，但它们的含义是相同的。

表17-1 Oracle与PostgreSQL数据类型对照表

数据类型	Oracle 11g	PostgreSQL 10
字符类型	VARCHAR、VARCHAR2、N VARCHAR、N VARCHAR2	character varying 或 text
	CHAR、NCHAR	character
	CLOB、NCLOB、LONG	text
数字类型	NUMBER	numeric
	FLOAT	real、double precision
(续)		
数据类型	Oracle 11g	PostgreSQL 10
时间日期类型	DATE	date 或 timestamp
	TIMESTAMP	timestamp
二进制类型	BLOB、RAW	bytea

Oracle 数据库迁移到 PostgreSQL 的 2 种方法

17-1 PostgreSQL 数据库迁移到 Oracle 数据库

Oracle 数据库迁移到 PostgreSQL 数据库

2. 数据库迁移

数据库迁移是指将数据从一个数据库迁移到另一个数据库的过程。数据库迁移可以分为物理迁移和逻辑迁移两种。物理迁移是指将数据库的物理文件（如数据文件、控制文件、日志文件等）从一个数据库迁移到另一个数据库。逻辑迁移是指将数据库的逻辑结构（如表、视图、索引等）从一个数据库迁移到另一个数据库。数据库迁移可以用于多种场景，如数据库升级、数据库迁移、数据库备份等。

PostgreSQL 数据库迁移到 Oracle 数据库的方法如下：
<https://www.postgresql.org/docs/10/static/pgsql.html>

17.3 数据库入门

数据库入门数据库入门SQL数据库入门数据库入门数据库入门SQL数据库入门SQL数据库入门SQL数据库入门数据库入门

数据库入门SQL数据库入门SELECT数据库入门UPDATE数据库入门INSERT数据库入门DELETE数据库入门SQL数据库入门SQL数据库入门数据库入门数据库入门Oracle数据库入门数据库入门数据库入门数据库入门SQL数据库入门数据库入门数据库入门数据库入门

1.SQL数据库入门

数据库入门SQL数据库入门PostgreSQL数据库入门Oracle数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门SQL数据库入门数据库入门

Oracle数据库入门ROWNUM数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门数据库入门

```
SQL> SELECT OBJECT_NAME FROM dba_objects WHERE ROWNUM < 2;
OBJECT_NAME
-----
-----
ICOL$
```

PostgreSQL LIMIT

```
postgres=# SELECT relname FROM pg_class LIMIT 1;
 relname 
-----
 test_sr
(1 row)
```

Oracle ROWNUM PostgreSQL LIMIT
ROWNUM
LIMIT

Oracle ROWNUM PostgreSQL LIMIT
ROWNUM

SQL Oracle
PostgreSQL Oracle
SQL

```
SQL> SELECT SEQ_1.CURRVAL FROM DUAL;
CURRVAL
-----
 2
```

PostgreSQL

```
postgres=# SELECT currval('seq_1');
currval
-----
2
(1 row)
```

currval() and nextval() are used to get the current and next value of a sequence. PostgreSQL and Oracle both support sequences.

Oracle and PostgreSQL both support sequences. Oracle sequences are created using the CREATE SEQUENCE statement, while PostgreSQL sequences are created using the CREATE SEQUENCE statement.

```
SQL> SELECT * FROM (SELECT * FROM table_1);
```

PostgreSQL sequences are created using the CREATE SEQUENCE statement.

```
mydb=> SELECT * FROM (SELECT * FROM table_1) as b;
```

SQL sequences are created using the CREATE SEQUENCE statement. Oracle sequences are created using the CREATE SEQUENCE statement. START WITH...CONNECT BY is used to create a sequence. PostgreSQL sequences are created using the CREATE SEQUENCE statement.

Oracle sequences are created using the CREATE SEQUENCE statement. t_area is a table. Sequences are used to generate unique values.

```
CREATE TABLE t_area(id numeric,name varchar(32),fatherid
numeric);
```

```
INSERT INTO t_area VALUES (1, '00' ,0);
INSERT INTO t_area VALUES (2, '00' ,1);
INSERT INTO t_area VALUES (3, '00' ,1);
INSERT INTO t_area VALUES (4, '00' ,2);
INSERT INTO t_area VALUES (5, '00' ,2);
INSERT INTO t_area VALUES (6, '00' ,3);
INSERT INTO t_area VALUES (7, '000' ,4);
INSERT INTO t_area VALUES (8, '000' ,4);
```

00ID00400000000000000000

```
SQL> SELECT id, name, fatherid
2 FROM t_area
3 START WITH id = 4
4 CONNECT BY PRIOR id = fatherid;
```

id	name	fatherid
4	00	2
7	000	4
8	000	4

PostgreSQL0000000000WITH
RECURSIVE000PostgreSQL000t_area0000
00000000000000

```
CREATE TABLE t_area(id int4,name varchar(32),fatherid int4);
```

00id0fatherid00000Oracle00numeric0
000PostgreSQL0000int40000000000000000

PostgreSQL mydb

```
mydb=> WITH RECURSIVE r AS (  
        SELECT * FROM t_area WHERE id = 4  
        UNION ALL  
        SELECT t_area.* FROM t_area, r WHERE  
t_area.fatherid = r.id  
    )  
    SELECT * FROM r ORDER BY id;  
 id | name | fatherid  
----+-----+-----  
  4 |   |      2  
  7 |   |      4  
  8 |   |      4  
(3 rows)
```

PostgreSQL 4.1

Oracle PostgreSQL SQL
Oracle PostgreSQL

2.

Oracle SQL Oracle
SQL PostgreSQL

Oracle SYSDATE CURRENT_DATE

```
SQL> SELECT SYSDATE,CURRENT_DATE FROM DUAL;
SYSDATE          CURRENT_DATE
-----
13-11-17         13-11-17
```

CURRENT_TIMESTAMP

```
SQL> SELECT CURRENT_TIMESTAMP FROM DUAL;
CURRENT_TIMESTAMP
-----
13-11-17 02.45.42.330637 +08:00
```

PostgreSQL CURRENT_DATE CURRENT_TIMESTAMP CURRENT_DATE CURRENT_TIMESTAMP now

```
mydb=> SELECT CURRENT_DATE;
current_date
-----
2017-11-12
(1 row)
```

```
mydb=> SELECT CURRENT_TIMESTAMP,now();
current_timestamp | now
-----+-----
2017-11-12 11:13:37.430499+08 | 2017-11-12
11:13:37.430499+08
(1 row)
```

Oracle INSTR
0 INSTR

{INSTR} (string , substring [, position [, occurrence]])

- string
 - substring
 - position1
 - occurrence1
-

```
SQL> SELECT INSTR('Hello PostgreSQL!','o') FROM DUAL;  
INSTR('HELLOPOSTGRESQL!','O')  
-----  
5
```

PostgreSQL

position(substring in string)

PostgreSQLのposition関数とOracleのINSTR関数の違い

```
mydb=> SELECT position('o' in 'Hello PostgreSQL!');
position
-----
5
(1 row)
```

PostgreSQLのposition関数は、文字列中の指定文字の位置を返す。OracleのINSTR関数は、文字列中の指定文字の位置を返す。ただし、OracleのINSTR関数は、文字列中の指定文字の位置を返す。ただし、OracleのINSTR関数は、文字列中の指定文字の位置を返す。

PostgreSQLのposition関数は、文字列中の指定文字の位置を返す。OracleのINSTR関数は、文字列中の指定文字の位置を返す。ただし、OracleのINSTR関数は、文字列中の指定文字の位置を返す。

17.4 数据迁移

本章主要介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库中。

- 如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库中，包括 text、csv 等格式。

- 如何在 PostgreSQL 中创建 oracle_fdw 扩展，并配置 PostgreSQL 与 Oracle 数据库的连接。

- 如何使用 ETL 工具进行数据迁移。

本章主要介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库中，包括 oracle_fdw 扩展的配置。

17-2 oracle_fdw 配置

角 色	主机名	IP	端 口	库 名	用户名	版 本	字符集
本地库 PostgreSQL	pghost1	192.168.28.74	1921	mydb	pguser	PostgreSQL10	UTF8
远程库 Oracle	pghost3	192.168.28.76	1521	oradb	community	Oracle 11g	UTF8

1.pgghost1 连接到 Oracle 11g

在pghost1上安装oracle_fdw包
Oracle包包括basic和devel
sqlplus包
RPM包
<http://www.oracle.com/technetwork/topics/inuxx86-64soft-092277.html>

安装RPM包

```
[root@pghost1 soft_bak]# rpm -ivh oracle-instantclient11.2-  
basic-11.2.0.1.0-1.x86_64.rpm  
[root@pghost1 soft_bak]# rpm -ivh oracle-instantclient11.2-  
sqlplus-11.2.0.1.0-1.x86_64.rpm  
[root@pghost1 soft_bak]# rpm -ivh oracle-instantclient11.2-  
devel-11.2.0.1.0-1.x86_64.rpm
```

安装postgres包
Oracle包

```
...PostgreSQL包  
export ORACLE_BASE=/usr/lib/oracle  
export ORACLE_HOME=/usr/lib/oracle/11.2/client64  
export  
LD_LIBRARY_PATH=$ORACLE_HOME/lib:$PGHOME/lib:/lib64:/usr/lib  
64:/usr/local/lib64:/lib:/usr/lib:/usr/local/lib  
export PATH=$ORACLE_HOME/bin:$PATH:$PGHOME/bin:.
```

安装Oracle包

在pghost1上安装postgres包
在pghost3上安装oracle包

```
[postgres@pghost1 ~]$ sqlplus
community/community@//192.168.28.76/oradb
SQL*Plus: Release 11.2.0.1.0 Production on Tue Nov 14
20:55:56 2017
Copyright (c) 1982, 2009, Oracle. All rights reserved.
Connected to:
Oracle Database 11g Enterprise Edition Release 11.2.0.1.0 -
64bit Production
With the Partitioning, OLAP, Data Mining and Real
Application Testing options

SQL>
```

安装Oracle数据库

2.pgghost1安装oracle_fdw

从https://api.pgxn.org/dist/oracle_fdw 安装
oracle_fdw数据库扩展
2.0.0

安装oracle_fdw

```
# unzip oracle_fdw-2.0.0.zip
```

安装PostgreSQL
root用户安装postgres数据库
扩展

```
[root@pghost1 oracle_fdw-2.0.0]# source
/home/postgres/.bash_profile
[root@pghost1 oracle_fdw-2.0.0]# which pg_config
/opt/pgsql/bin/pg_config
```

编译安装oracle_fdw

```
[root@pghost1 oracle_fdw-2.0.0]# make
...
[root@pghost1 oracle_fdw-2.0.0]# make install
/bin/mkdir -p '/opt/pgsql_10.0/lib'
/bin/mkdir -p '/opt/pgsql_10.0/share/extension'
/bin/mkdir -p '/opt/pgsql_10.0/share/extension'
/bin/mkdir -p '/opt/pgsql_10.0/share/doc/extension'
/usr/bin/install -c -m 755 oracle_fdw.so
'/opt/pgsql_10.0/lib/oracle_fdw.so'
/usr/bin/install -c -m 644 ./oracle_fdw.control
'/opt/pgsql_10.0/share/extension/'
/usr/bin/install -c -m 644 ./oracle_fdw--1.1.sql
./oracle_fdw--1.0--1.1.sql
'/opt/pgsql_10.0/share/extension/'
/usr/bin/install -c -m 644 ./README.oracle_fdw
'/opt/pgsql_10.0/share/doc/extension/'
```

安装后的文件位置

`$PGHOME/share/extension`

```
[root@pghost1 oracle_fdw-2.0.0]# ll
/opt/pgsql/share/extension/oracle_fdw*
-rw-r--r-- 1 root root 231 Nov 12 14:45
/opt/pgsql/share/extension/oracle_fdw--1.0--1.1.sql
-rw-r--r-- 1 root root 1003 Nov 12 14:45
/opt/pgsql/share/extension/oracle_fdw--1.1.sql
-rw-r--r-- 1 root root 133 Nov 12 14:45
/opt/pgsql/share/extension/oracle_fdw.control
```

```
sudo su - postgres
cd $PGHOME/share/extension
make
make install
sudo su - postgres
create extension oracle_fdw;
```

3.pgghost1 mydb oracle_fdw

```
pgghost1 mydb oracle_fdw
```

```
[postgres@pgghost1 ~]$ psql mydb postgres
psql (10.0)
Type "help" for help.
```

```
mydb=# CREATE EXTENSION oracle_fdw;
CREATE EXTENSION
```

```
CREATE FOREIGN DATA WRAPPER oracle_fdw
    FROM 'oracle_fdw'
    HANDLER 'pguser'
    VALIDATOR 'pguser';
```

```
mydb=# GRANT USAGE ON FOREIGN DATA WRAPPER oracle_fdw TO
pguser;
GRANT
```

4.pgghost3 Oracle

```
community oracle
community
```

```
CREATE TABLE T_ORA1(id numeric,info varchar2(32),create_time
timestamp);
```

```
INSERT INTO T_ORA1 VALUES (1,'a',CURRENT_TIMESTAMP);
INSERT INTO T_ORA1 VALUES (2,'b',CURRENT_TIMESTAMP);
INSERT INTO T_ORA1 VALUES (3,'',CURRENT_TIMESTAMP);
```

Oracle
PostgreSQL
Oracle

```
CREATE USER READONLY IDENTIFIED BY "readonly"
  DEFAULT TABLESPACE TS_COMMUNITY
  TEMPORARY TABLESPACE TEMP
  PROFILE DEFAULT
  ACCOUNT UNLOCK;
```

```
GRANT CONNECT TO READONLY;
GRANT SELECT ON COMMUNITY.T_ORA1 TO READONLY;
```

T_ORA1
READONLY

5.pgghost1mydb

pgusermydb

```
[postgres@pgghost1 ~]$ psql mydb pguser
psql (10.0)
Type "help" for help.
```

```
mydb=> CREATE SERVER fs_oracle_pgghost3 FOREIGN DATA WRAPPER
```

```
oracle_fdw OPTIONS (dbserver '//192.168.28.76/oradb');  
CREATE SERVER
```

```
  OPTIONS (dbserver '//192.168.28.76/oradb');  
CREATE SERVER oracle_fdw FOREIGN DATA WRAPPER Oracle  
OPTIONS (dbserver '//192.168.28.76/oradb');  
CREATE SERVER oracle_fdw FOREIGN DATA WRAPPER Oracle  
OPTIONS (dbserver '//192.168.28.76/oradb');
```

```
CREATE SERVER oracle_fdw FOREIGN DATA WRAPPER Oracle  
OPTIONS (dbserver '//192.168.28.76/oradb');
```

```
oradb=  
  (DESCRIPTION=  
    (ADDRESS=  
      (PROTOCOL=TCP)  
      (HOST=192.168.28.76)  
      (PORT=1521)  
    )  
    (CONNECT_DATA=  
      (SERVICE_NAME=oradb)  
    )  
  )  
)
```

```
CREATE SERVER oracle_fdw FOREIGN DATA WRAPPER Oracle  
OPTIONS (dbserver '//192.168.28.76/oradb');
```

```
mydb=> CREATE SERVER fs_oracle_pghost3 FOREIGN DATA WRAPPER  
oracle_fdw OPTIONS (dbserver 'oradb');
```

dbserver
db

```
mydb=> CREATE USER MAPPING FOR pguser SERVER
fs_oracle_pghost3 OPTIONS (user 'readonly', password
'readonly');
CREATE USER MAPPING
```

pguser PostgreSQL
OPTIONS readonly Oracle

6.pghost1mydb

mydb Oracle T_ORA1
T_ORA1

```
mydb=> CREATE FOREIGN TABLE ft_t_ora1 (
    id    int4,
    info  text,
    create_time timestamp with time zone
) SERVER fs_oracle_pghost3 OPTIONS (schema 'COMMUNITY',
table 'T_ORA1');
```

SERVER
fs_oracle_pghost3

OPTIONS

·Table Oracle Oracle
Oracle
Oracle

·Schema Oracle

·Readonly yes
INSERT UPDATE DELETE
false

schema table
https://pgxn.org/dist/oracle_fdw/

```
mydb=> SELECT * FROM ft_t_oral;
   id | info | create_time
-----+-----+-----
   1 | a    | 2017-11-12 15:02:54.465304+08
   2 | b    | 2017-11-12 15:02:54.477453+08
   3 |     | 2017-11-12 15:02:54.493603+08
(3 rows)
```

PostgreSQL Oracle

7. Oracle PostgreSQL

PostgreSQL Oracle


```
mydb=> CREATE TABLE t_oral(id int4,info text,create_time
timestamp with time zone);
CREATE TABLE
```

INSERT

```
mydb=> INSERT INTO t_oral SELECT * FROM ft_t_oral;  
INSERT 0 3
```

[illegible]

```
mydb=> SELECT * FROM t_oral;
```

id	info	create_time
1	a	2017-11-12 15:02:54.465304+08

```

      2 | b          | 2017-11-12 15:02:54.477453+08
      3 |          | 2017-11-12 15:02:54.493603+08
(3 rows)

```

OracleT_OA1
PostgreSQL

17.5 □□□□□□□□

SQL DBA
DBA
SQL

SQL DBA

17.6 数据库

数据库是存储和管理数据的系统。数据库系统由数据库、数据库管理系统、数据库管理员、数据库用户等组成。数据库系统的主要功能是数据的存储、检索、更新、删除、插入、备份、恢复等。

数据库系统的主要组成部分包括：数据库、数据库管理系统、数据库管理员、数据库用户。数据库系统是数据库管理员使用数据库管理系统对数据库进行管理的系统。

数据库系统的主要组成部分包括：数据库、数据库管理系统、数据库管理员、数据库用户。数据库系统是数据库管理员使用数据库管理系统对数据库进行管理的系统。数据库系统的主要组成部分包括：数据库、数据库管理系统、数据库管理员、数据库用户。

17.7 oracle_fdw

oracle_fdw

oracle_fdw.so

oracle_fdw

```
mydb=# CREATE EXTENSION oracle_fdw;
ERROR:  could not load library
"/opt/pgsql_10.0/lib/oracle_fdw.so": libclntsh.so.11.1:
cannot open shared object file: No such file or directory
```

```
lib libOracle
Oracle
$LD_LIBRARY_PATH $ORACLE_HOME
lib
$PGHOME/lib
```

```
[root@pghost1 ~]# cp
/usr/lib/oracle/11.2/client64/lib/libclntsh.so.11.1
/opt/pgsql/lib/
[root@pghost1 ~]# cp
/usr/lib/oracle/11.2/client64/lib/libnnz11.so
/opt/pgsql/lib/
```

oracle_fdw

Oracle 11gR2 64bit 安装

Oracle 11gR2 64bit 安装

```
mydb=> SELECT * FROM ft_t_oral;  
ERROR:  error connecting to Oracle: OCIEnvCreate failed to  
create environment handle  
DETAIL:
```

Oracle 11gR2 64bit 安装
\$ORACLE_HOME \$LD_LIBRARY_PATH
Oracle 11gR2 64bit 安装

```
... PostgreSQL  
export ORACLE_BASE=/usr/lib/oracle  
export ORACLE_HOME=/usr/lib/oracle/11.2/client64  
export  
LD_LIBRARY_PATH=$ORACLE_HOME/lib:$PGHOME/lib:/lib64:/usr/lib  
64:/usr/local/lib64:/lib:/usr/lib:/usr/local/lib  
export PATH=$ORACLE_HOME/bin:$PATH:$PGHOME/bin:.
```

Oracle 11gR2 64bit 安装
sqlplus
Oracle 11gR2 64bit 安装

```
sqlplus community/community@//192.168.28.76/oradb
```

Oracle 11gR2 64bit 安装
sqlplus Oracle 11gR2 64bit 安装
PostgreSQL 10.4 安装

创建表

创建表 oracle_fdw 用于连接 Oracle 数据库
Oracle 与 PostgreSQL 的连接
PostgreSQL 连接

创建表 Oracle 连接

创建表 oracle_fdw 用于连接 Oracle 数据库
table 用于连接 Oracle 数据库
创建表

创建表 ft_t_ora1 用于连接 Oracle 数据库
创建表

```
CREATE FOREIGN TABLE ft_t_ora2 (  
    id      int4,  
    info    text,  
    create_time timestamp with time zone  
) SERVER fs_oracle_pghost3 OPTIONS (schema 'COMMUNITY',  
table 't_ora1');
```

创建表 ft_t_ora2 用于连接 Oracle 数据库

```
mydb=> SELECT * FROM ft_t_ora2;  
ERROR:  Oracle table "COMMUNITY"."t_ora1" for foreign table  
"ft_t_ora2" does not exist or does not allow read access  
DETAIL:  ORA-00942: table or view does not exist  
HINT:  Oracle table names are case sensitive (normally all  
uppercase)
```

ft_t_ora2 Oracle
ft_t_ora2

schema
ft_t_ora3 ft_t_ora1 schema
ft_t_ora3

```
CREATE FOREIGN TABLE ft_t_ora3 (  
    id      int4,  
    info    text,  
    create_time timestamp with time zone  
) SERVER fs_oracle_pghost3 OPTIONS (schema 'community',  
table 'T_ORA1');
```

ft_t_ora3

```
mydb=> SELECT * FROM ft_t_ora3;  
ERROR:  Oracle table "community"."T_ORA1" for foreign table  
"ft_t_ora3" does not exist or does not allow read access  
DETAIL:  ORA-00942: table or view does not exist  
HINT:   Oracle table names are case sensitive (normally all  
uppercase).
```

ft_t_ora3

17.8 数据迁移

数据迁移是指将数据从一个数据库迁移到另一个数据库。Oracle 和 PostgreSQL 都是常用的数据库。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

Oracle 数据库和 PostgreSQL 数据库都是常用的数据库。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

<http://ora2pg.darold.net/> 是一个开源的 Oracle 数据库到 PostgreSQL 数据库的数据迁移工具。本文将介绍如何将 Oracle 数据库中的数据迁移到 PostgreSQL 数据库。

18 PostGIS

PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。

PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。PostgreSQL 数据库系统是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。

PostGIS 是一个开源的、关系型的数据库系统，它支持多种数据类型，包括文本、数字、日期、时间、几何类型等。

18.1 安装

安装PostGIS扩展包

使用yum或apt-get或brew安装PostgreSQL和PostGIS扩展包

```
[postgres@pghost1 ~]$ yum search postgres
Loaded plugins: fastestmirror, security
...
...
...
postgis24_10.x86_64 : Geographic Information Systems
Extensions to PostgreSQL
    Name and summary matches only, use "search all" for
    everything.
```

安装PostgreSQL 10和PostGIS扩展包

```
[postgres@pghost1 ~]$ yum install -y postgis24_10
```

安装PostGIS扩展包



18.2 安装GIS扩展

PostGIS是PostgreSQL的一个扩展，它提供了空间数据的支持。安装PostGIS的步骤如下：

```
postgres=# CREATE DATABASE geo;
CREATE DATABASE
postgres=# \c geo
geo=# CREATE EXTENSION postgis;
CREATE EXTENSION
```

安装PostgreSQL时，PostGIS扩展可能不会自动安装。您可以通过以下命令检查已安装的扩展：

```
geo=# SELECT name,default_version FROM
pg_available_extensions WHERE name ~ 'gis';
      name      | default_version
-----+-----
 postgis        | 2.4.3
 postgis_tiger_geocoder | 2.4.3
 postgis_topology | 2.4.3
 postgis_sfcgal  | 2.4.3
```

您可以通过以下命令查看PostGIS的完整版本信息：

```
gis=# SELECT postgis_full_version();
POSTGIS="2.4.3 r16312" PGSQL="100" GEOS="3.6.2-CAPI-1.10.2
4d2925d6" PROJ="Rel. 4.9.3, 15 August 2016" GDAL="GDAL
```

1.11.5, released 2016/07/01" LIBXML="2.9.7" LIBJSON="0.12.1"
RASTER

GIS

18.3 空間DB

PostGISは空間データを扱うためのPostgreSQLの拡張機能です。



PostGISはPostgreSQLの空間データを扱うための拡張機能です。PostGISはSpatial Typeと呼ばれる空間データを扱うための型を提供します。

PostGISはPostgreSQLの空間データを扱うための拡張機能です。Geometryと呼ばれる空間データを扱うための型を提供します。PostGISはgeosと呼ばれる空間データを扱うためのライブラリを提供します。

18.3.1 空間データの形式

- PostGISは空間データを扱うための拡張機能です。
- ・空間データを扱うためのWKT(Well-Known-Text)
 - ・空間データを扱うためのWKB(Well-Known-Binary)
 - ・GeoJSON

12

18.3.2 ☐ ☐ ☐ ☐ ☐ ☐

```

    Point ST_Point
  }
};

```

```
geo=# SELECT Point(1,2), ST_Point(1,2);
 point |          st_point
-----+-----
  (1,2) | 010100000000000000000000F03F0000000000000040
```

PostgreSQLのPointはDoubleを
16ビットでPostGISはST_Pointは
21ビットで表現する。これは
ID

PostGISはPostgreSQLの
ETL
pg_dump

ST_AsTextはWKT
WKB

```
geo=# SELECT ST_AsText(ST_Point(1,2));
 st_astext
-----
POINT(1 2)
```

PostGISは
Geometry
Geometry

```
gis=# CREATE TABLE geo (
    geom GEOMETRY
);
gis=# INSERT INTO geo VALUES
    (ST_Point(1.0, 2.0)),
    ('LineString(0 0,1 1,2 1,2 3)'),
    ('Polygon((0 0, 1 0, 1 1,0 1,0 0))'),
    ('MultiPoint(1 2,3 4)');
```

18.3.3 空间数据类型

空间数据类型包括点、线、面、多面体、几何集合等。在PostgreSQL中，空间数据类型可以通过WKT或GeoJSON格式进行存储。

```
gis=# SELECT
      ST_AsText(geom)    AS wkt,
      ST_AsGeoJSON(geom) AS json
FROM geo;
```

图18-1 空间数据类型

图18-1 空间数据类型

wkt	json
POINT(1 2)	{"type":"Point","coordinates":[1,2]}
LINESTRING(0 0,1 1,2 1,2 3)	{"type":"LineString","coordinates":[[0,0],[1,1],[2,1],[2,3]]}
POLYGON((0 0,1 0,1 1,0 1,0 0))	{"type":"Polygon","coordinates":[[[0,0],[1,0],[1,1],[0,1],[0,0]]]}
MULTIPOINT(1 2,3 4)	{"type":"MultiPoint","coordinates":[[1,2],[3,4]]}

空间数据类型在PostgreSQL中可以通过ST_AsText和ST_AsGeoJSON函数进行转换。

```
geo=# SELECT
      ST_AsLatLonText(ST_Point(116.321367,39.966956));
      st_aslatlon
-----
39°58'1.042"N 116°19'16.921"E
```

18.3.4 空間データベース

PostGISは空間データベースとして、PostgreSQLに空間データを扱うための機能を追加したデータベースです。

PostGISは、PostgreSQLの空間データを扱うための機能を追加したデータベースです。

```
geo=# SELECT ST_Point(1,1) <-> ST_Point(2,2);
1.4142135623730951
```

PostGISは、PostgreSQLの空間データを扱うための機能を追加したデータベースです。

PostGISは、PostgreSQLの空間データを扱うための機能を追加したデータベースです。

```
gis=# SELECT ST_Point(116.321367, 39.966956) <->
ST_Point(116.315346, 39.997398);
-- 0.03103172255933419
```

4326 WGS84
GPS 3.4km

```
gis=# SELECT ST_AsText(ST_GeomFromText('POINT(116.321367
39.966956)', 4326)) :: GEOGRAPHY <->
      ST_AsText(ST_GeomFromText('POINT(116.315346 39.997398)',
4326)) :: GEOGRAPHY;
3423.653480690467
```

R ST_length
WGS84

```
-- R MultiLineString
gis=# SELECT
ST_length(ST_GeomFromText('MULTILINESTRING((116.351494
39.976407,116.353159 39.976395,116.353365 39.976479),
(116.350712 39.976471,116.350984 39.976406,116.351494
39.976407,116.351755 39.976479),(116.351516
39.976558,116.346836 39.976502),(116.350984
39.976406,116.345421 39.976355,116.344646
39.976304,116.343376 39.976326,116.343353
39.976185,116.343285 39.976147,116.340918 39.976076),
(116.332306 39.976515,116.326279 39.976308),(116.354314
39.976488,116.347592 39.976436,116.347855 39.97637),
(116.354302 39.976574,116.351516 39.976558,116.351407
39.976616,116.349808 39.976603,116.346841
39.976569,116.346836 39.976502,116.34048
39.976443,116.335918 39.976538),(116.343194
39.976399,116.340565 39.976374,116.333921
39.976471,116.331997 39.976417,116.332003 39.976342),
(116.346841 39.976569,116.344304 39.97658,116.34256
39.976526,116.336051 39.976581,116.335918
39.976538,116.333922 39.976571,116.332306
39.976515,116.332108 39.976573,116.326431
39.976373,116.326018 39.976391),(116.331997
39.976417,116.323548 39.976145),(116.347592
```

```
39.976436,116.343194 39.976399,116.343376
39.976326,116.340209 39.976302,116.336219
39.976364,116.335423 39.976342,116.333617
39.976397,116.327658 39.976201,116.326968
39.976147,116.325406 39.976134,116.324303
39.97606,116.323696 39.976076),(116.317584
39.97613,116.320373 39.976186,116.323807
39.976325,116.324613 39.976315,116.325891
39.976391,116.326018 39.976391,116.326279
39.976308,116.317613 39.976048),(116.31983
39.975958,116.32005 39.976033,116.317835
39.975977,116.317637 39.975969,116.317777
39.97589,116.323696 39.976076,116.323548 39.976145,116.32005
39.976033))',4326)::geography);
12314.809569165007
```

□□□□□□□□R□□□□□□□□□□□□12□□□□

□□□□□□□□□□□□□□□□□□□□ST_Area□
ST_Polygon□□□□□□□□□□□□

```
gis=# SELECT ST_Area(ST_GeomFromText('POLYGON((116.402408
39.915326,116.402822 39.91533,116.402847
39.915036,116.402434 39.915035,116.402426
39.91495,116.402474 39.91397,116.402493 39.91367,116.402505
39.913385,116.40251 39.913269,116.401675 39.91323,116.400835
39.913238,116.400346 39.913231,116.398979
39.913198,116.398359 39.913159,116.398337
39.913153,116.398352 39.912786,116.398317
39.912738,116.398274 39.9127,116.398221 39.912661,116.39806
39.912645,116.398029 39.912642,116.39792
39.912639,116.397909 39.912638,116.397904 39.912668,116.3971
39.912644,116.396687 39.91261,116.396659
39.912609,116.396553 39.912606,116.39653
39.912605,116.396519 39.912746,116.396401
39.912747,116.396354 39.912759,116.396331 39.91294,116.39632
39.913079,116.396302 39.913093,116.396264
39.913105,116.396167 39.913097,116.394376
39.91302,116.393154 39.912964,116.392198
39.912919,116.392031 39.912918,116.392004
```


18.4 空间数据库的创建

首先LBS系统需要建立一个空间数据库，用于存储位置信息。这里我们使用PostgreSQL 9.1.3版本，并安装PostGIS 2.0.1版本。

在Ubuntu 12.04 LTS上，安装PostGIS的步骤如下：
1. 安装PostgreSQL 9.1.3：
sudo apt-get install postgresql-9.1.3
2. 安装PostGIS 2.0.1：
sudo apt-get install postgresql-postgis-2.0.1

```
gis=# CREATE TABLE stations(  
    name      TEXT,  
    longitude DOUBLE PRECISION,  
    latitude  DOUBLE PRECISION,  
);  
-- 创建索引  
gis=# SELECT name FROM stations WHERE longitude BETWEEN  
116.312358 AND 116.330376  
AND latitude BETWEEN 39.957947 AND 39.975965;
```

接下来，我们需要创建一个GeoHash索引，用于快速查找位置信息。这里我们使用PostGIS 2.0.1中的GeoHash函数。

最后，我们使用PostGIS 2.0.1中的SQL函数来查询位置信息。

```
gis=# SELECT
    name,
    ST_Point(116.321367, 39.966956) :: GEOGRAPHY <->
position :: GEOGRAPHY AS distance
FROM stations
WHERE ST_Point(116.321367, 39.966956) :: GEOGRAPHY <->
position :: GEOGRAPHY < 500
ORDER BY ST_Point(116.321367, 39.966956) :: GEOGRAPHY <->
position :: GEOGRAPHY;
```

18.4.1 空间索引

在100米范围内查找所有POI点，并按距离排序。A在1000米范围内。

```
gis=# SELECT name FROM poi ORDER BY position <->
ST_Point(116.458855, 39.909863) LIMIT 1000;
```

在PostgreSQL中使用GIST索引。

```
gis=# CREATE INDEX CONCURRENTLY idx_poi_position_gist ON poi
USING gist (position);
```

QUERY PLAN

```

-----
Limit (cost=0.42..9.73 rows=10 width=31)
  -> Index Scan using idx_poi_position_gist on poi
(cost=0.42..58440964.86 rows=62750132 width=31)
    Order By: ("position" <->
'0101000000CAA65CE15D1D5D40946B0A6476F44340'::geometry)

```

□□□□□□□1ms□□□□□□□□□□□□□□□□

```

geo=# SELECT name FROM poi ORDER BY position <->
ST_Point(116.458855, 39.909863) LIMIT 10;
      name
-----
...
...
...
(10 rows)
Time: 0.993 ms

```

18.4.2 □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□

```

-- □□□□ (AOI, Area of Interest)
gis=# CREATE TABLE aoi(
      name    TEXT,
      bound   GEOMETRY

```



```
)  
--  A  
gis=# SELECT name FROM aoi WHERE ST_Contains(bound,  
ST_Point(116.458855, 39.909863));  
...  
...  
...  
...  
...
```

18.5 数据库

PostGIS数据库是PostgreSQL数据库的一个扩展，它提供了对地理数据的存储、查询和管理的功能。PostGIS数据库是PostgreSQL数据库的一个扩展，它提供了对地理数据的存储、查询和管理的功能。

PostGIS数据库是PostgreSQL数据库的一个扩展，它提供了对地理数据的存储、查询和管理的功能。PostGIS数据库是PostgreSQL数据库的一个扩展，它提供了对地理数据的存储、查询和管理的功能。